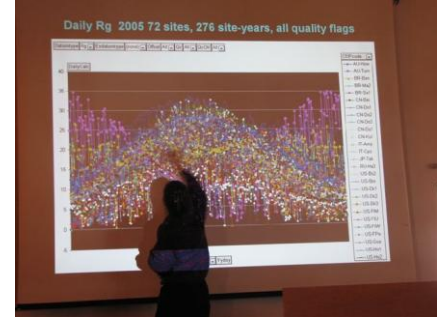


Life in the Cloud: MODISAzure

Catharine van Ingen
Partner Architect
eScience Group, Microsoft Research

Introduction



- ▶ In the summer of 2009, I started MODIS Azure with three goals:
 - Learn more about remote sensing, maps, and geospatial data
 - Kick some Azure tires
 - Do some science I could understand
- ▶ Just over a year later, we're doing amazing science and computer science
 - All of that is thanks to terrific science collaborators and computer science collaborators
 - The challenges ahead are new and definitely 4th paradigm
- ▶ This talk highlights specific learnings and potential opportunities for collaboration forward

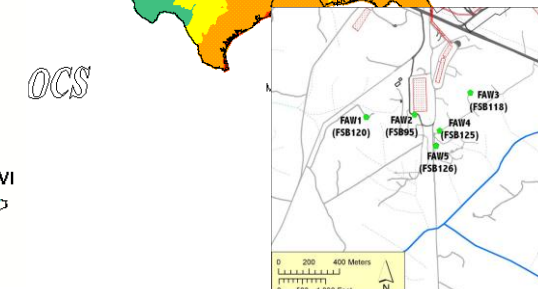
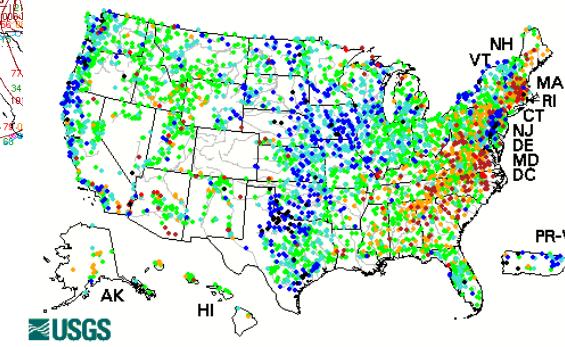
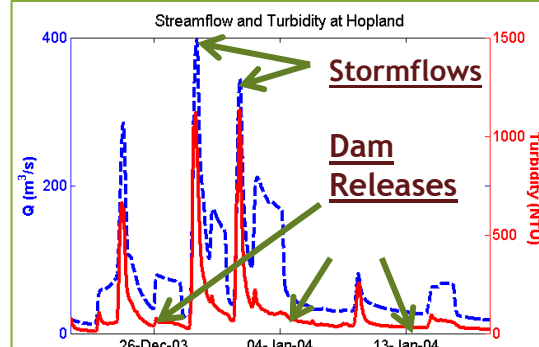
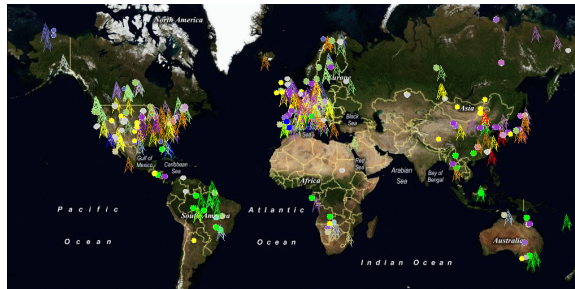
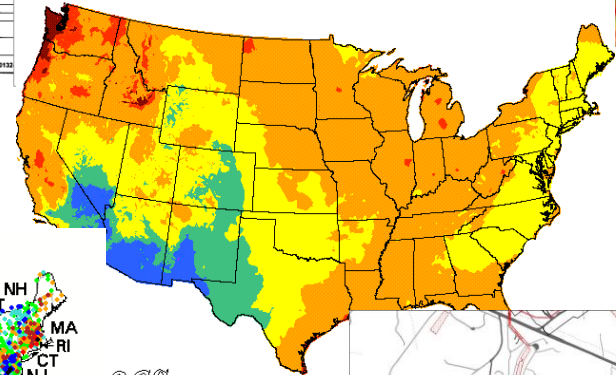
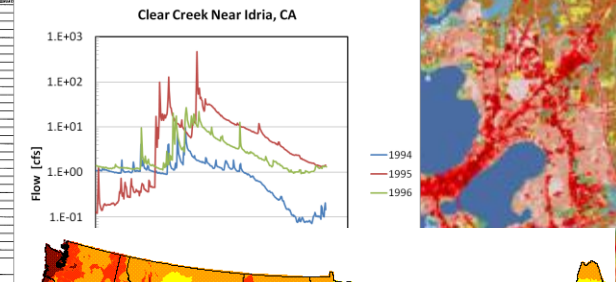
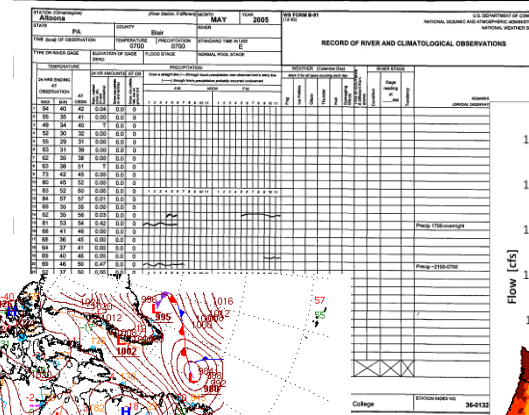
Environmental Science Today

Big Science ! Hallelujah!
Big Science ! Yodelie Hoo!
Laurie Anderson

The Data (and Truth?) Are Out There

- ▶ National and International Datasets
 - USGS National Water Information System
 - NOAA National Climatic Data Center
 - FLUXNET Network
 - Satellite data (e.g. MODIS)
- ▶ Local Datasets
 - Local Agencies
 - Companies (e.g. Timber)
 - Ecology Organizations
 - Individual Researchers

Page	Target	Habitat Attribute	Indicator	Method	Status	Poor	Fair
1	Spawning Adult	Estuary	Passage at Mouth	Pool Option		<30 days	>30 days
2	Spawning Adult	Hydrology	Passage Point	Flow Panel Results	Done		
3	Spawning Adult	Passage	Physical Barriers	Passage Database	Final	<50% of IP-km	50-70%
4	Spawning Adult	Viability	Freshwater Harvest	Review Regulations	Status?		
5	Spawning Adult	Viability	Density Target	NMFS Calculation	Apply TRT Criteria	Watershed Specific	
6	Spawning Adult	Sediment	Spawning Gravel	Take all talus with emb. rating <5, multiply by avg width of riffle squared	Hopland Doing Queries		
7							



Data Variety – The Spice of Life



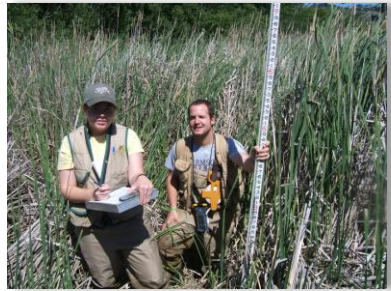
Manual Measurement



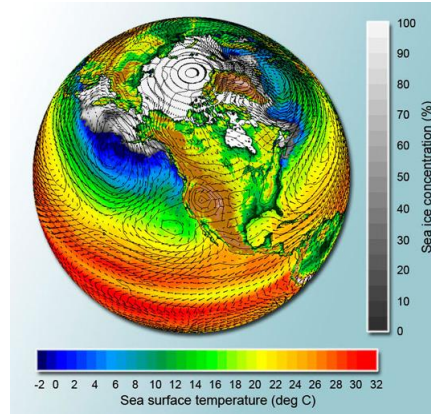
Automated Measurement



Sample Collection



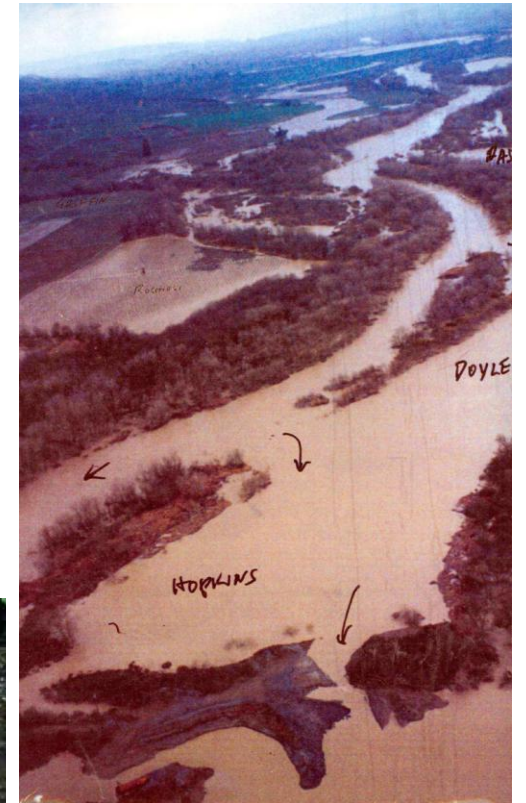
Typing



Model Output



Counting



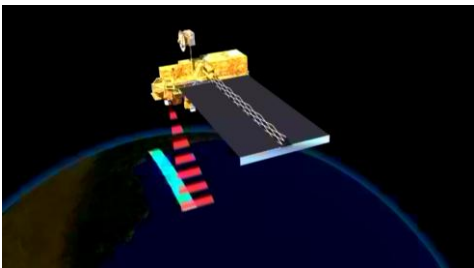
Historical Photographs



Relatively Ubiquitous Motes

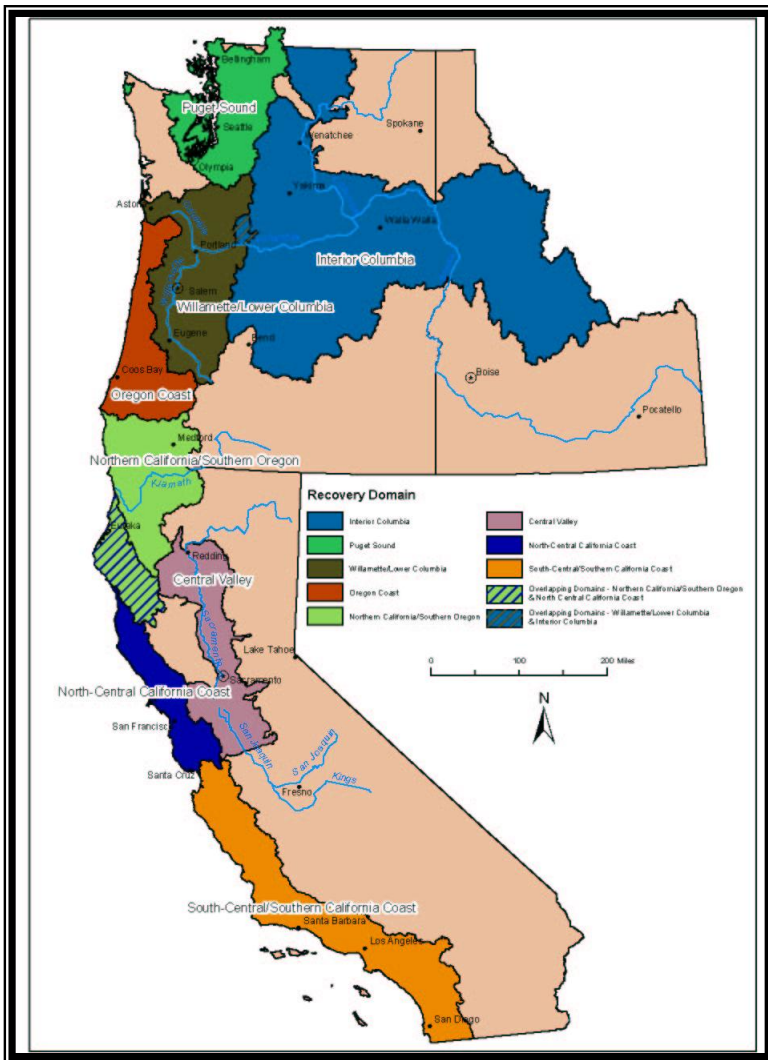


Aircraft Surveys

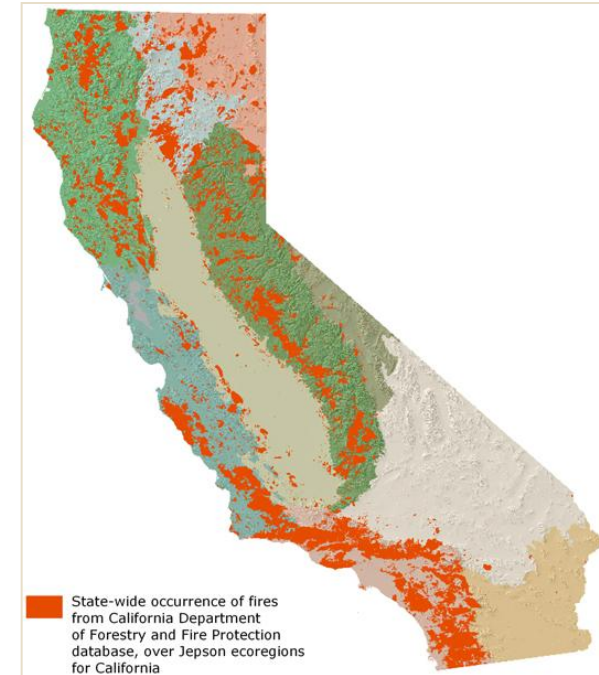


Satellite

Global and Regional Eco-Science



- ▶ A major shift is happening in the way eco-science is done.
 - Moving from individual studies of local processes to collaborative studies of regional and global processes. (e.g. studying the impact of climate change)
- ▶ Studying global scale environmental processes requires:
 - Integration of local, regional, and global spatial scales.
 - Integration across disciplines, e.g., climatology, hydrology, forestry, etc., and across methodologies (field observations, remote sensing, and modeling).



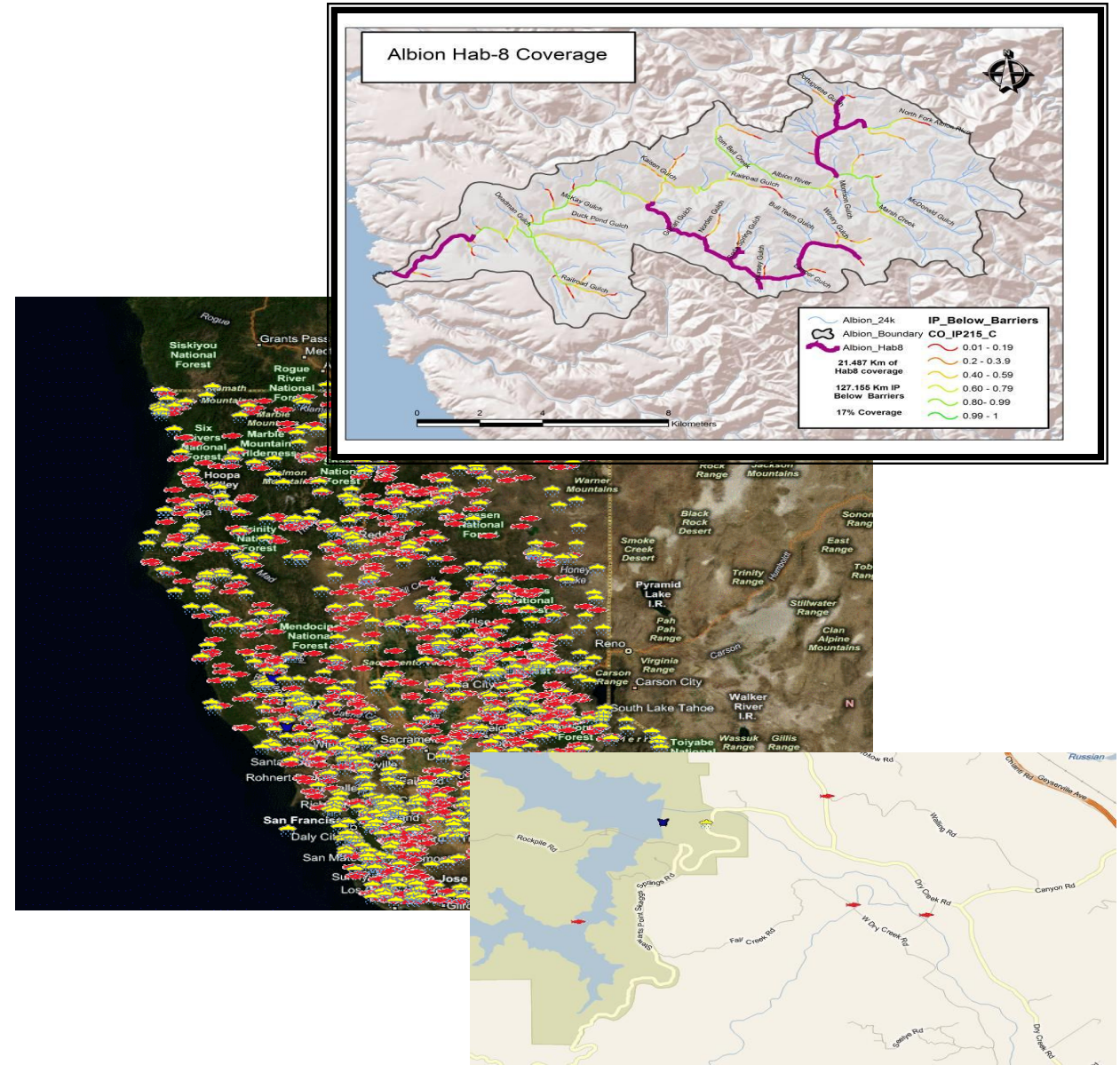
Synthesis –Bringing It All Together

USGS stream gauges
Coho presence/absence data
MODIS evapotranspiration
Cross section locations
Samples from 1997
Stream temperature gauges
Water system
Obstructions and other
Human activities

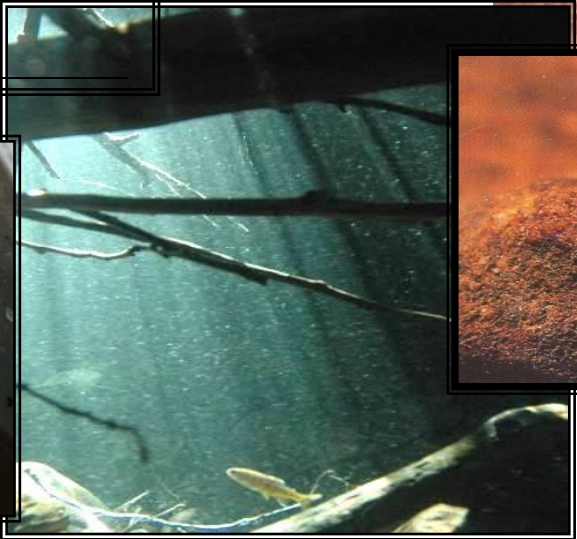
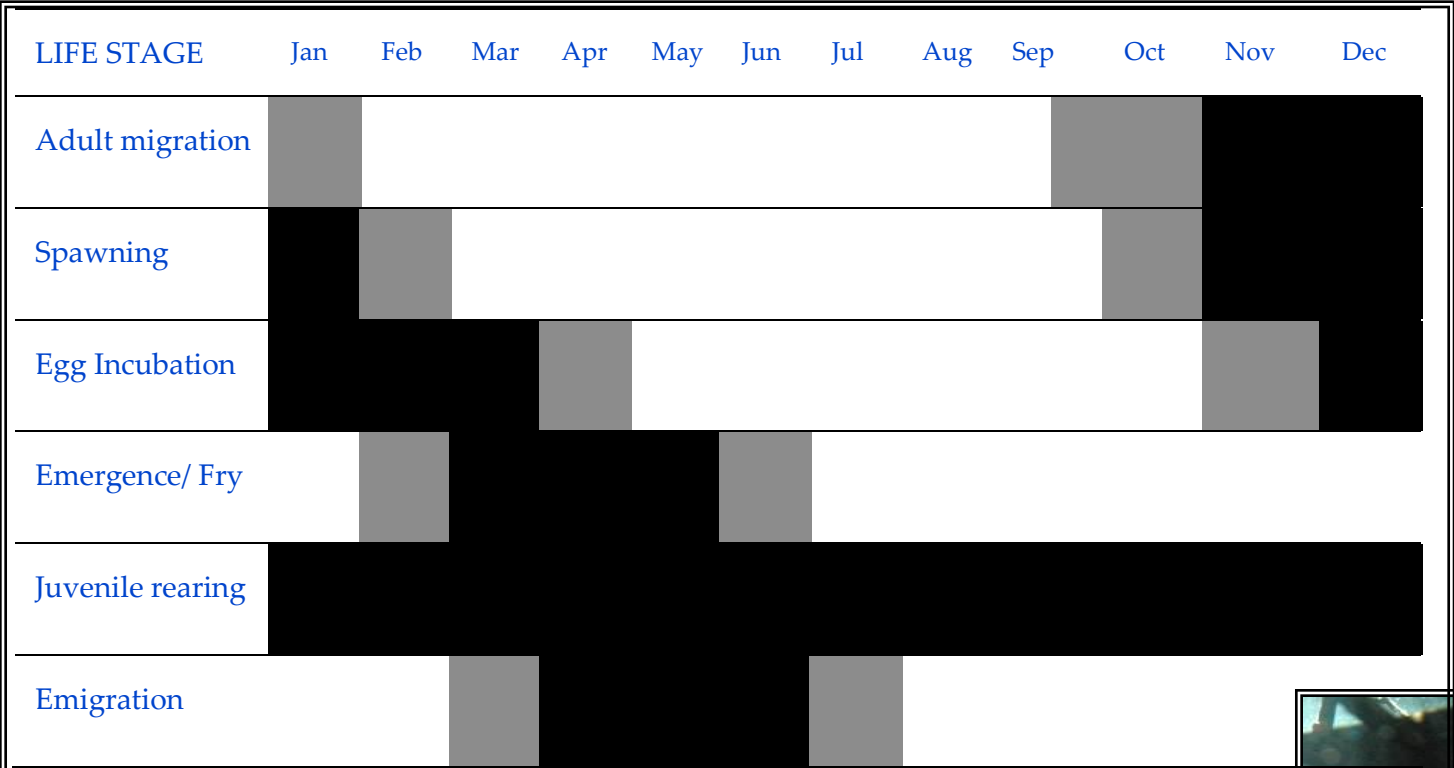


Data Integration Challenges

- ▶ Regular rasters, points, and spatial features
- ▶ Time series and intermittent
- ▶ Vocabulary meanings (ontology)
- ▶ Sparse in time, duration, or location
- ▶ Science variable derivation
- ▶ Gaps
- ▶ Spatial/temporal harmonization



Time is Not Just Another Axis : Salmon Year

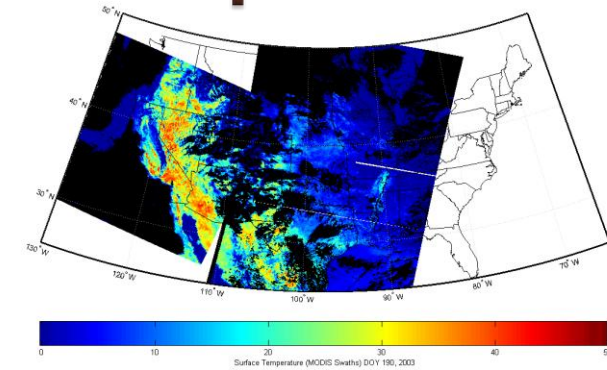


Tiling: Do Scientists have to become Computer Scientists?

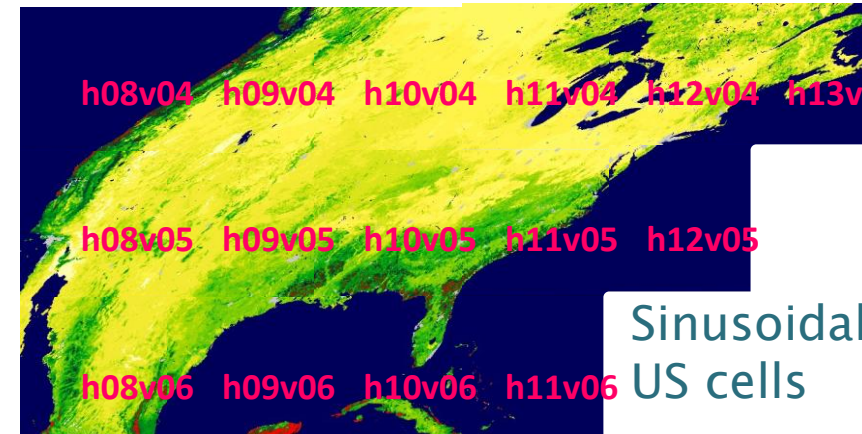
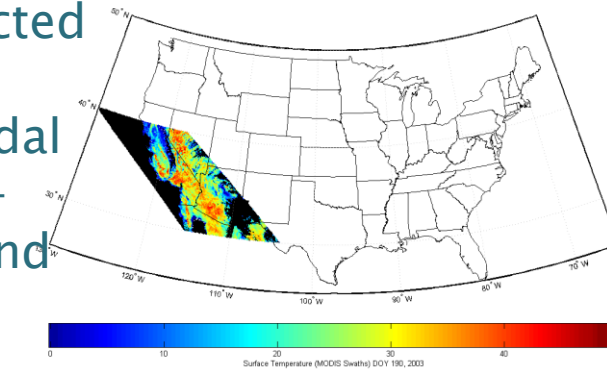
- ▶ Reprojection
 - Converts one geo-spatial representation to another.
 - Example: latitude-longitude swaths converted to sinusoidal cells.
- ▶ Spatial resampling
 - Converts one spatial resolution to another.
 - Example is converting from 1 KM to 5 KB pixels.
- ▶ Temporal resampling
 - Converts one temporal resolution to another.
 - Example is converting from daily observation to 8 day averages.
- ▶ Gap filling
 - Assigns values to pixels without data either due to inherent data issues such as clouds or missing pixels.
- ▶ Masking
 - Eliminates uninteresting or unneeded pixels.
 - Examples are eliminating pixels over the ocean when computing a land product or outside a spatial feature such as a watershed.

Grunge means you're doing science

Source
Data
(Swath
format)



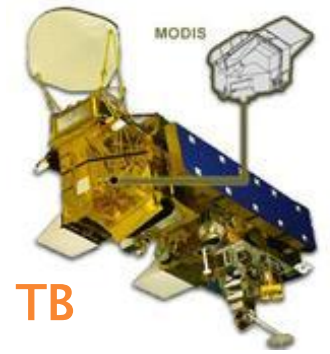
Reprojected
Data
(Sinusoidal
format –
equal land
area
pixel)



Sinusoidal
US cells

Wait ! It's Worse Than That!

- ▶ Provenance and trust widely varies
 - Data acquisition, early processing, and reporting ranges from a large government agency to individual scientists.
 - Smaller data often passed around in email; big data downloads can take days (if at all)
 - Opaque safe-deposit boxes and storage lockers prevail today
- ▶ Data sharing concerns and patterns vary
 - Open access followed by (non-repeatable and tedious) pre-processing
 - True science ready data set but concerns about misuse, misunderstanding particularly for hard won data.
- ▶ Computational tools differ.
 - Not everyone can get an account at a supercomputer center
 - Very large computations require engineering (error handling)
 - Space and time aren't always simple dimensions



Science happens when PBs, TBs, GBs, and KBs can be mashed up simply

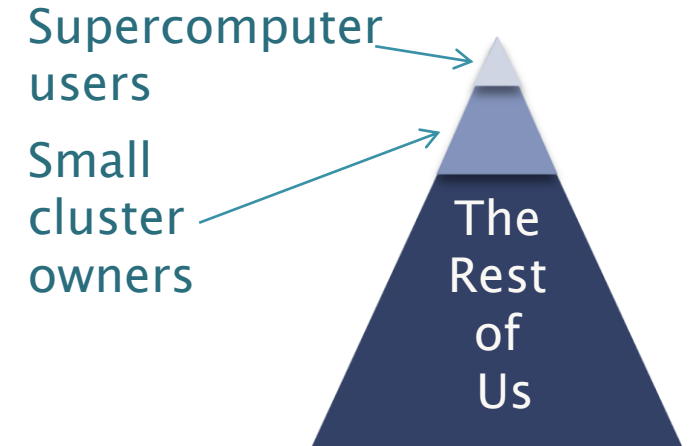
Bridging the Gap with the Cloud

► Barriers to Science:

- Resource: compute, storage, networking, visualization capability
- Complexity: specific cross-domain knowledge
- Tedium: repetitive data gathering or preprocessing tasks

► With cloud computing, we can:

- marshal needed storage and compute resources on demand without caring or knowing how that happens
- access living curated datasets without having to find, educate, and reward a private data curator
- run key common algorithms as Software as a Service without having to know the coding details or installing software
- grow a given collaboration or share data and algorithms across science collaborations elastically



Where do you want your data?



Democratizing science analysis by fostering sharing and reuse

Computing Evapotranspiration (ET)

You never miss the water till the well has run dry
Irish Proverb

What is ET ?

- ▶ Evapotranspiration (ET) is the release of water to the atmosphere by evaporation from open water bodies and transpiration, or evaporation through plant membranes, by plants.
- ▶ Climate change isn't just about a change in temperature, it's also about a change in the water balance and hence water supply critical to human activity.

From Dr. Youngryel Ryu's science research proposal:

Evapotranspiration (E) is a major component of the terrestrial hydrological cycle (ca. 60% of precipitation) [Trenberth, et al., 2007]. It controls land-atmosphere feedbacks and constitutes an important source of water vapor to the atmosphere [Raupach, 1998]. In turn, atmospheric water vapor is the most significant greenhouse gas and thus plays a fundamental role in weather and climate [Held and Soden, 2000]. Understanding E is important for socio-economic reasons, such as regulating available water for human use [Brauman, et al., 2007]. Thus, there have been diverse efforts to regularly monitor E in a regional scale using satellite remote sensing imagery [Anderson, et al., 2008; Diak, et al., 2004; Nishida, et al., 2003].

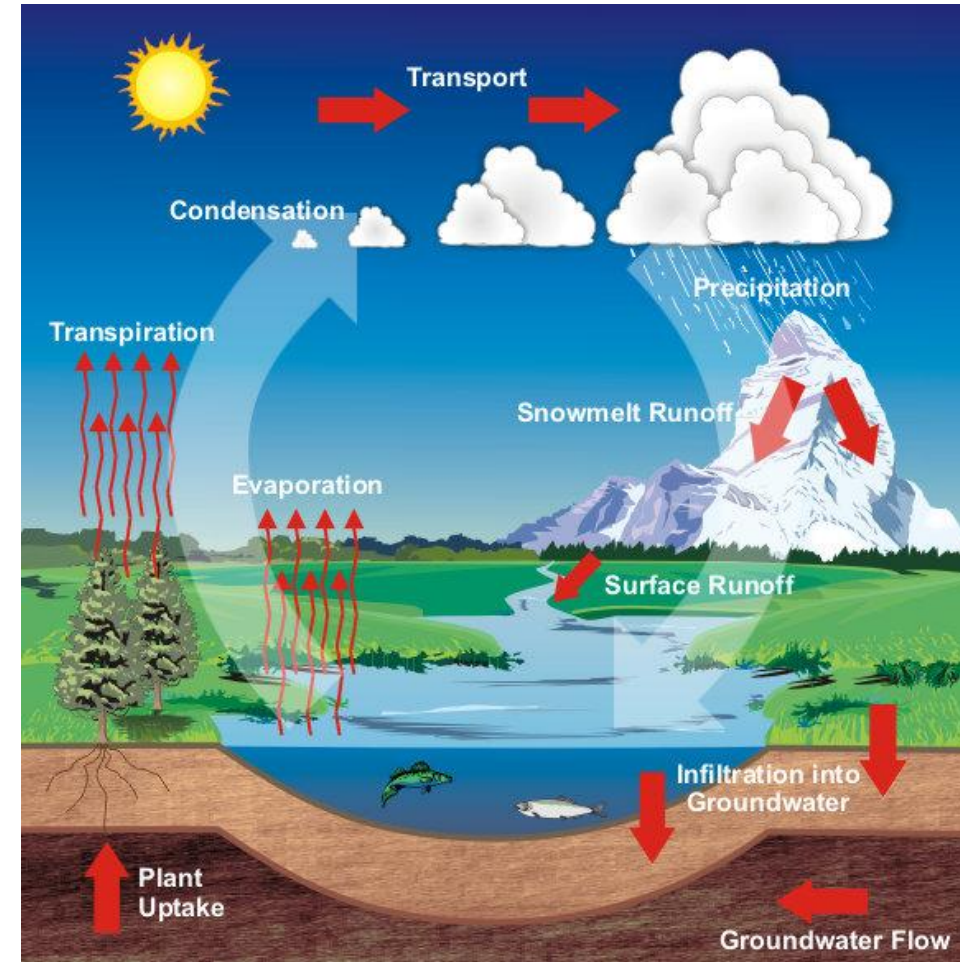


Image courtesy of the
[National Oceanic and Atmospheric Administration](#)

Computing ET From Historical Sensor Data

$$ET = P - R - \frac{dS}{dt}$$

Simple Water Balance

ET: Evapotranspiration or release of water to the atmosphere by evaporation from open water bodies and transpiration by plants

P: Precipitation including snowfall

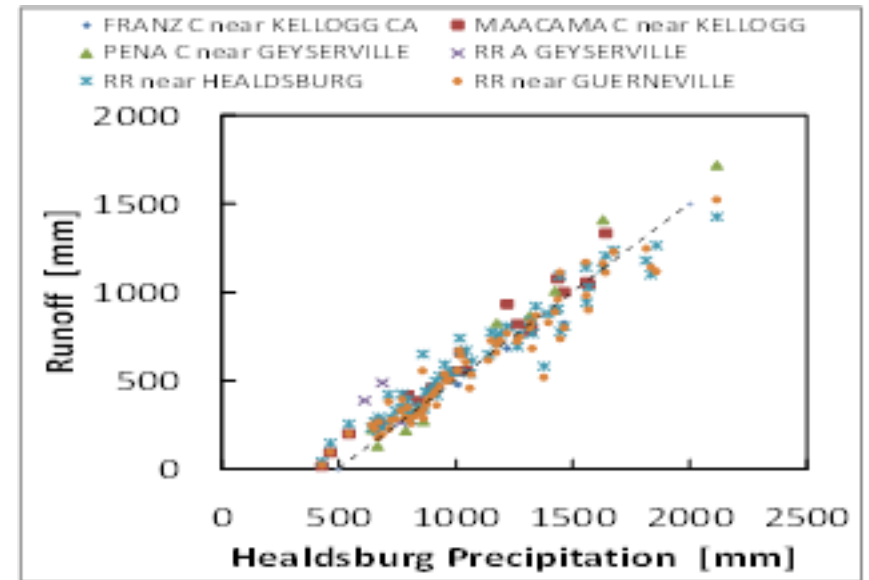
R: Surface runoff in streams and rivers

dS/dt : change in water storage over time such as increase in lakes or groundwater levels

P: <http://www.ncdc.noaa.gov/oa/ncdc.html>

R: <http://waterdata.usgs.gov/nwis>

- ▶ Easy to do (with a digital watershed)
- ▶ Long term trends only



In Mediterranean climates such as California, a long term equilibrium may exist. The ecosystem determines ET by soils and climate and the lowest recorded annual rainfall may determine vegetation.

~400 MB of data reduced to ~1 KB

Computing ET from First Principles

$$ET = \frac{\Delta R_n + \rho_a c_p (\delta q) g_a}{(\Delta + \gamma(1 + g_a/g_s)) \lambda_v}$$

Penman–Monteith (1964)

ET = Water volume evapotranspired ($\text{m}^3 \text{s}^{-1} \text{m}^{-2}$)

Δ = Rate of change of saturation specific humidity with air temp. (Pa K^{-1})

λ_v = Latent heat of vaporization (J/g)

R_n = Net radiation (W m^{-2})

c_p = Specific heat capacity of air ($\text{J kg}^{-1} \text{K}^{-1}$)

ρ_a = dry air density (kg m^{-3})

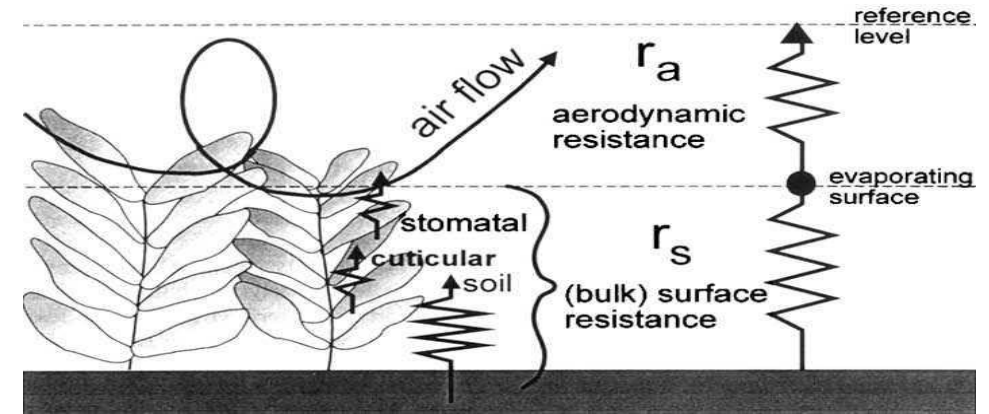
δq = vapor pressure deficit (Pa)

g_a = Conductivity of air (inverse of r_a) (m s^{-1})

g_s = Conductivity of plant stoma, air (inverse of r_s) (m s^{-1})

γ = Psychrometric constant ($\gamma \approx 66 \text{ Pa K}^{-1}$)

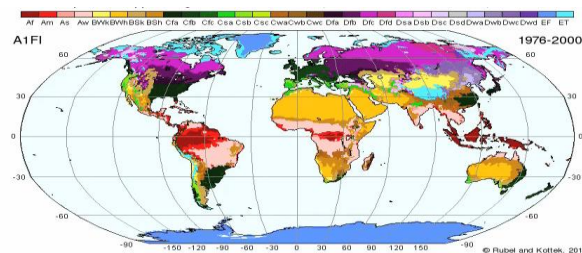
- ▶ Lots of inputs : big reduction
- ▶ Some of the inputs are not so simple
- ▶ Many have categorical dependencies



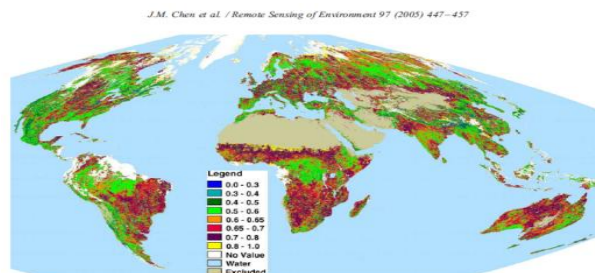
Estimating resistance/conductivity across a catchment can be tricky



Computing ET from Imagery, Sensors and Field Data

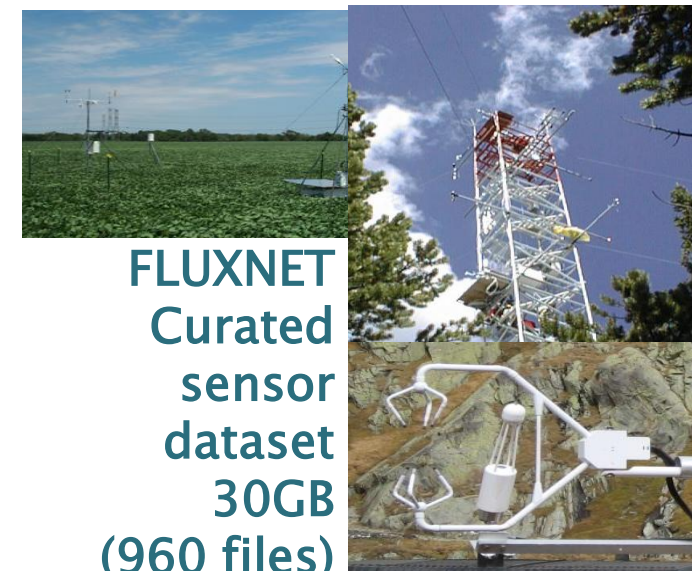
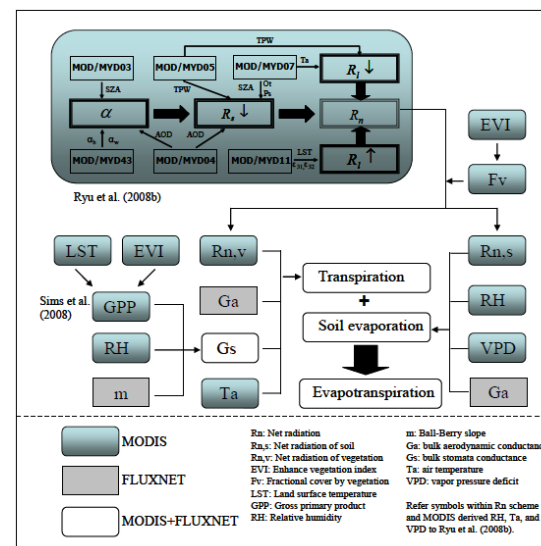


Climate classification
~1MB (1 file)

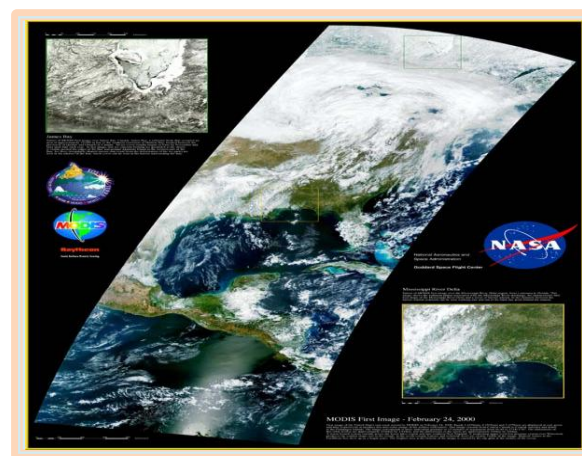


Vegetative clumping
~5MB (1 file)

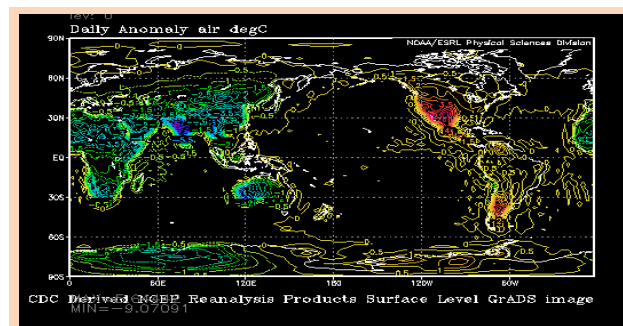
Not just a simple matrix computation due to dry region leaf/air temperatures differences, snow cover, leaf area fill, temporal upscaling, gap fill, biome conductance lookup, C3/C4 plants, etc etc



FLUXNET
Curated sensor dataset
30GB
(960 files)



NASA MODIS imagery archives
5 TB (600K files) for 10 US years



NCEP/NCAR ~100MB
(4K files)

FLUXNET
curated field dataset
2 KB (1 file)



MODISAzure : Computing ET in The Cloud

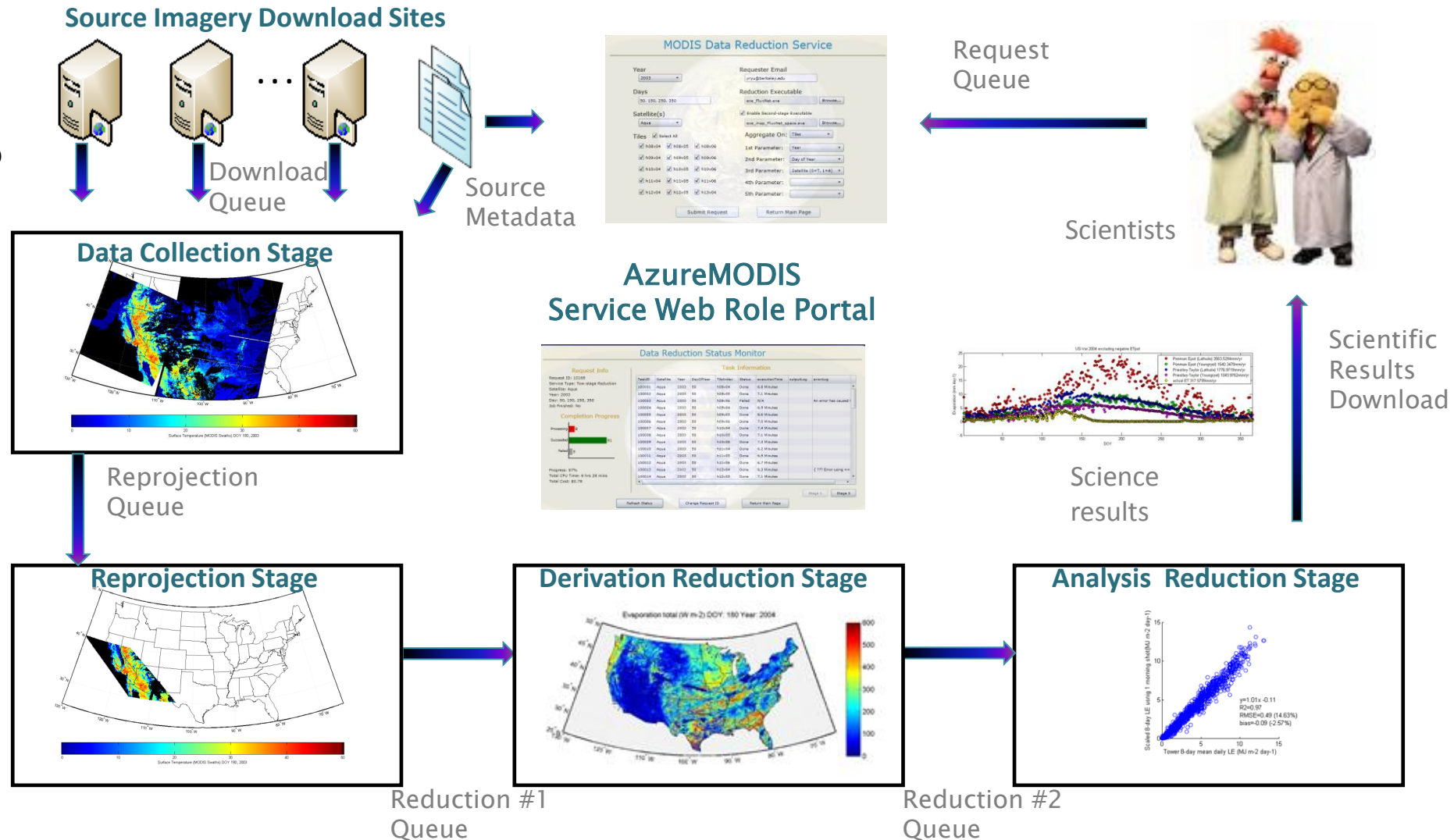
Behind every cloud is another cloud.

Judy Garland



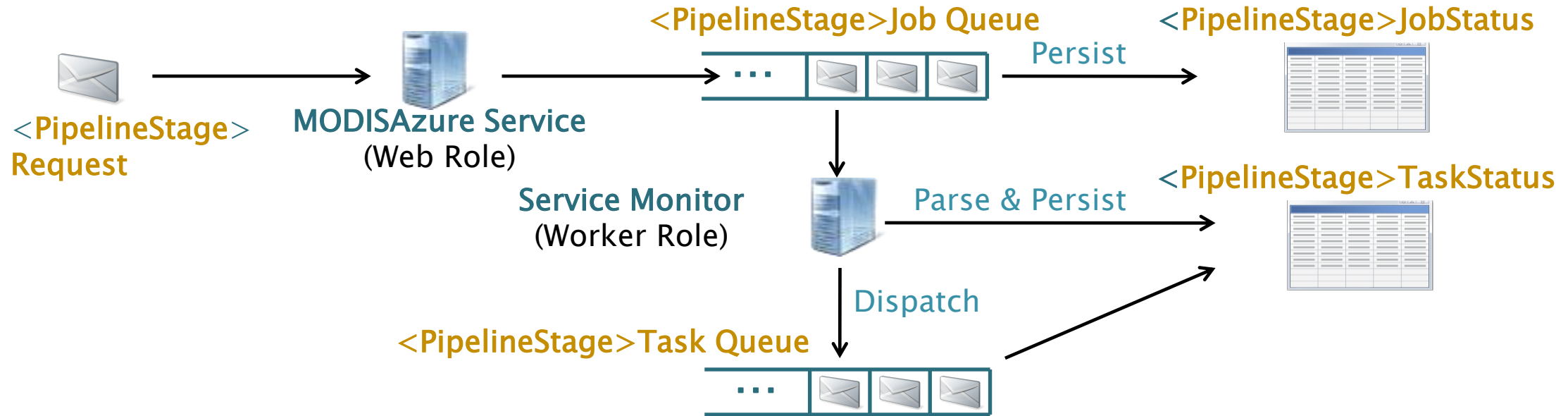
MODIS Azure: Four Stage Image Processing Pipeline

- ▶ Data collection stage
 - Downloads requested input tiles from NASA ftp sites
 - Includes geospatial lookup for non-sinusoidal tiles that will contribute to a reprojected sinusoidal tile
- ▶ Reprojection stage
 - Converts source tile(s) to intermediate result sinusoidal tiles
 - Simple nearest neighbor or spline algorithms
- ▶ Derivation reduction stage
 - First stage visible to scientist
 - Computes ET in our initial use
- ▶ Analysis reduction stage
 - Optional second stage visible to scientist
 - Enables production of science analysis artifacts such as maps, tables, virtual sensors



<http://research.microsoft.com/en-us/projects/azure/azuremodis.aspx>

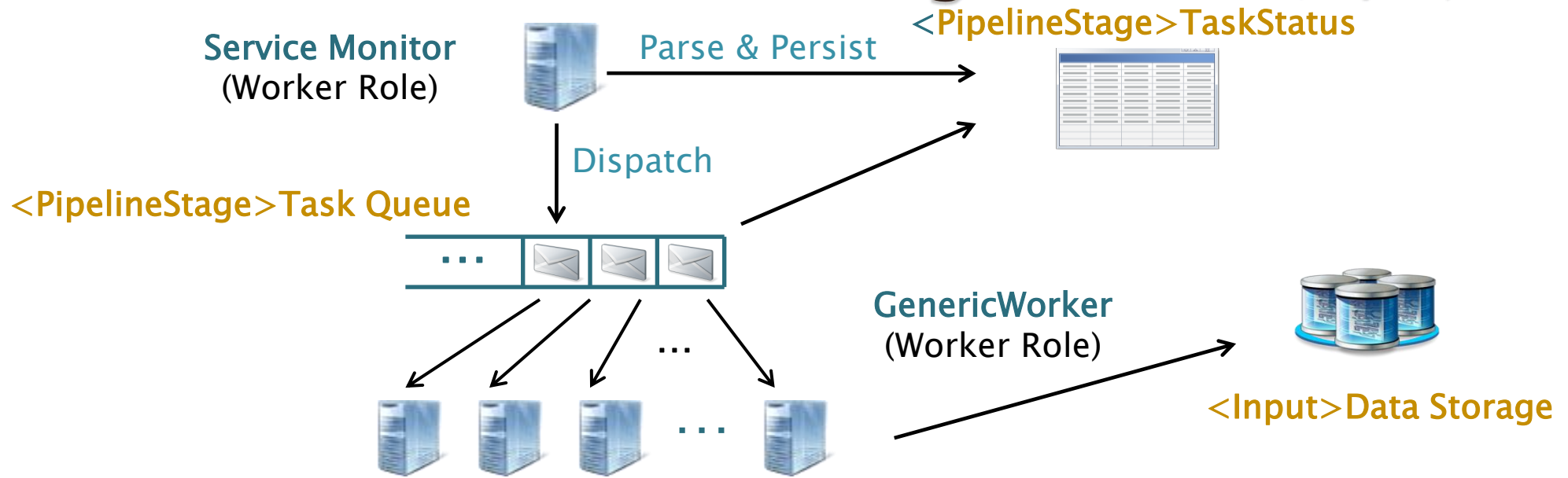
MODISAzure: Architectural Big Picture (1 / 2)



- ▶ **ModisAzure Service** is the Web Role front door
 - Receives all user requests
 - Queues request to appropriate Download, Reprojection, or Reduction Job Queue

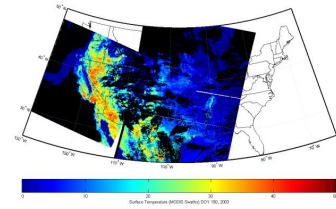
- ▶ **Service Monitor** is a dedicated Worker Role
 - Parses all job requests into tasks – recoverable units of work
 - Execution status of all jobs and tasks persisted in Tables

MODISAzure: Architectural Big Picture (2/2)



- ▶ All work actually done by a **GenericWorker** Worker Role
 - Dequeues tasks created by the Service Monitor
 - Science executable is sandboxed on an Azure Worker instance thereby enabling simple desktop development and debug
 - Marshalls all storage from/to Azure blob storage to/from local Azure Worker instance files
 - Retries failed tasks 3 times
 - Maintains all task status

Determining Inputs



- ▶ Each science variable is associated with a MODIS product
 - Terra satellite products (eg MOD04) used preferentially as they tend to be day time observations
 - Aqua satellite products (eg MYD04) used when Terra products unavailable
 - MCD products are higher level products generated by a combination of Terra and Aqua
 - MOD44B from 2000 used throughout
- ▶ Each product is either **swath** or **sinusoidal** projection
 - Sinusoidal are ready to use
 - Groups of swath products must be reprojected to create a sinusoidal tile
- ▶ Each product has a recurrence interval of daily, 8 day, 16 day, annual

M*D04	Aerosol
M*D05	Precipitable water
M*D06	Cloud
M*D07	Temperature, ozone
MCD43B*	Albedo
M*D11	Surface temperature
M*D15	LAI
MOD13A2	Vegetation Index
MCD12Q1	Land Cover
MOD44B	Veg. Contig. Fields

Determining What to Download

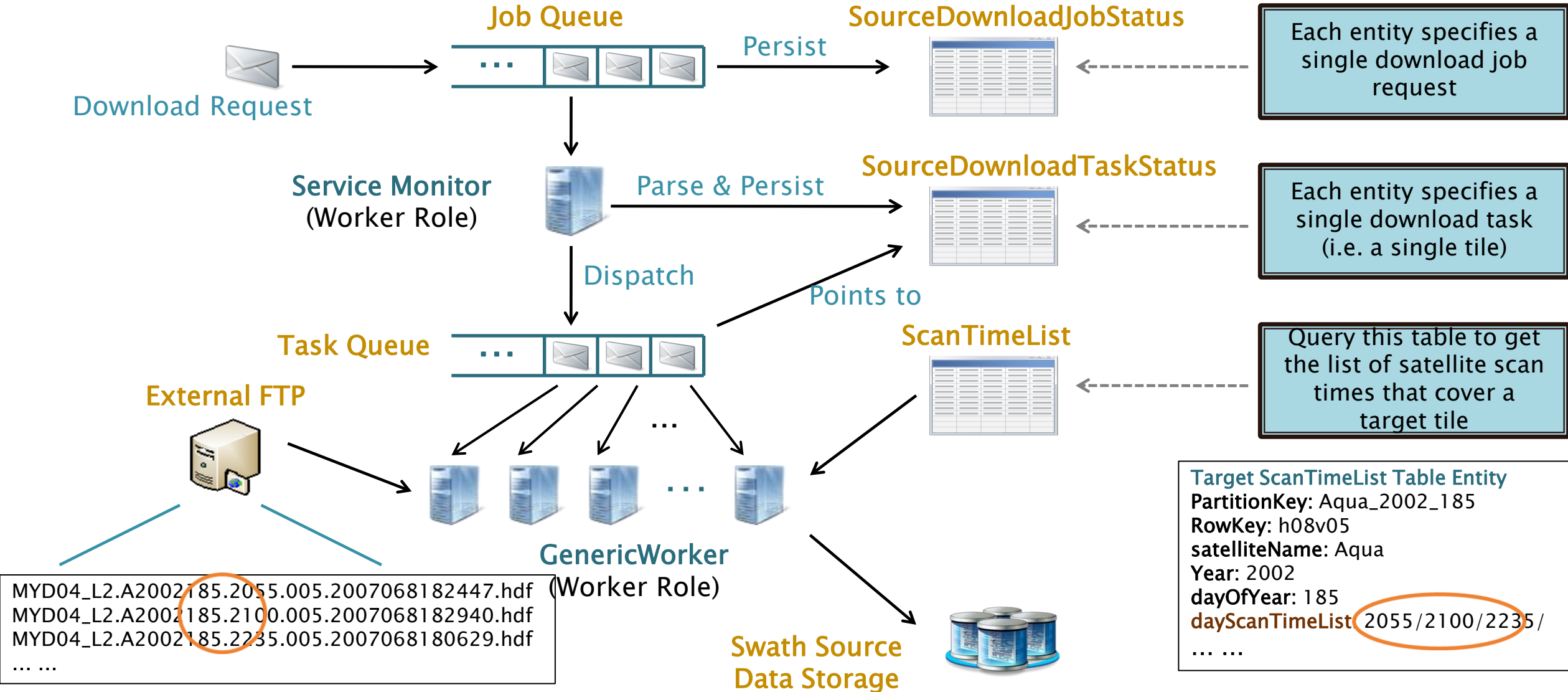
- ▶ NASA publishes a geometadata information for Terra and Aqua
- ▶ For each 5 minute swath data file (or granule) on the ftp site there is a corresponding geometa file containing:
 - DayNightFlag indicating day, night or both
 - Corner point latitude/longitude,
 - Bounding coordinates
- ▶ We ingested all files (288 per day * 10 years * 2 satellites) into a SQL database then paged the information into our Azure ScanTimeList Table
- ▶ The dayScanTimeList in the ScanTimeList table identifies all swath source file precursors for a given sinusoidal tile and drives the download and reprojection

#Attributes	PartitionK	RowKey	Timestamp	betweenScanTimeList	dayOfYear	dayScanTimeList	hIndex	nightScanTimeList	satellite	vIndex	year	
Terra	2003	160	h00v07	2/10/2010 7:33		160	2220/2355/	0	1005/1010/1145/	Terra	0	2003

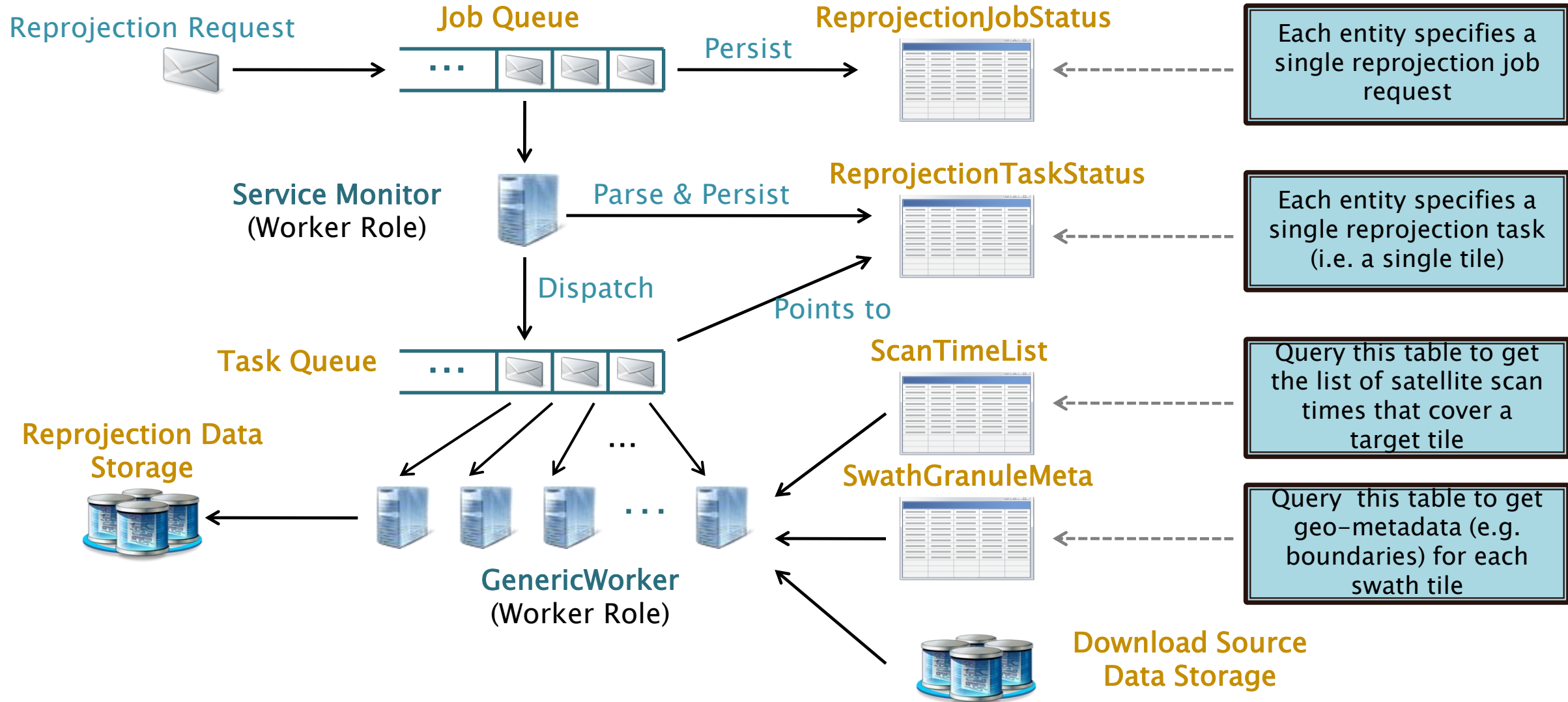
<ftp://ladsweb.nascom.nasa.gov/geoMeta/README>

Source Data Download Service

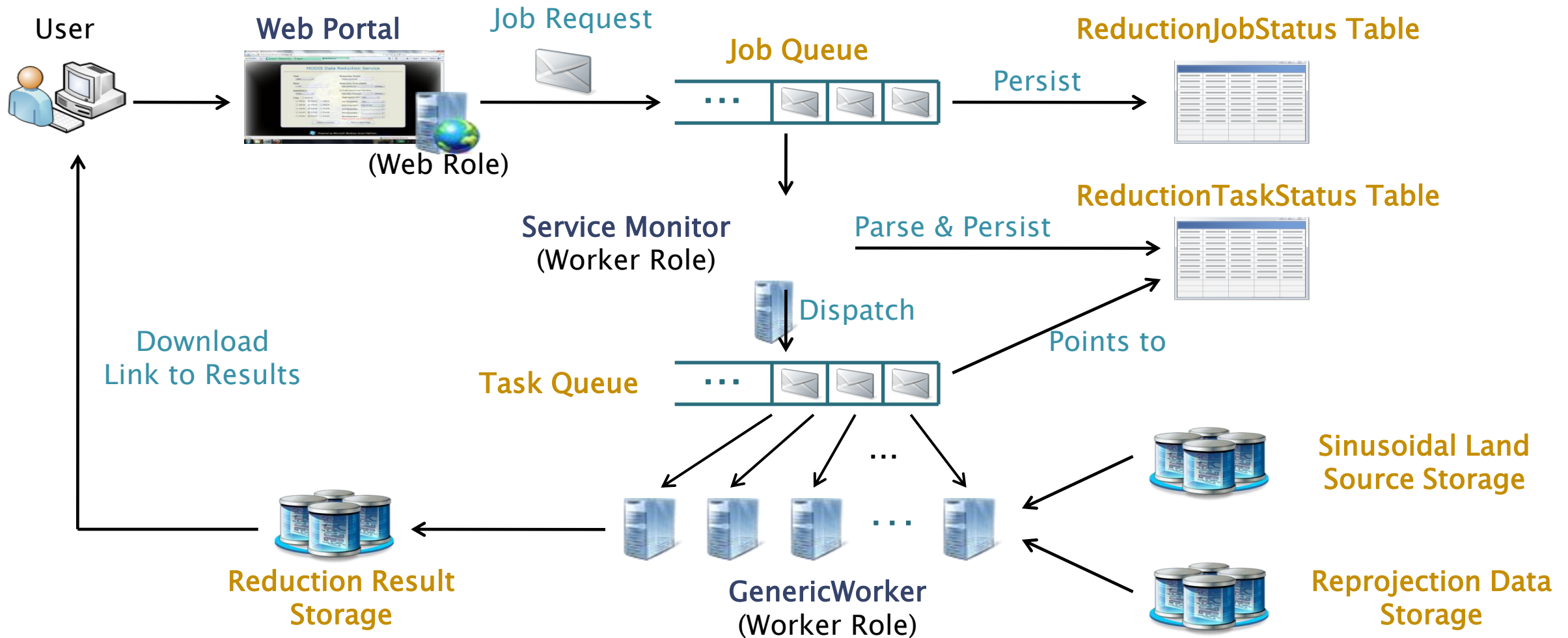
Example: Download the required source files for the target sinusoidal tile: MYD04_L2, Year 2002, Day 185, h08v05



Reprojection Service

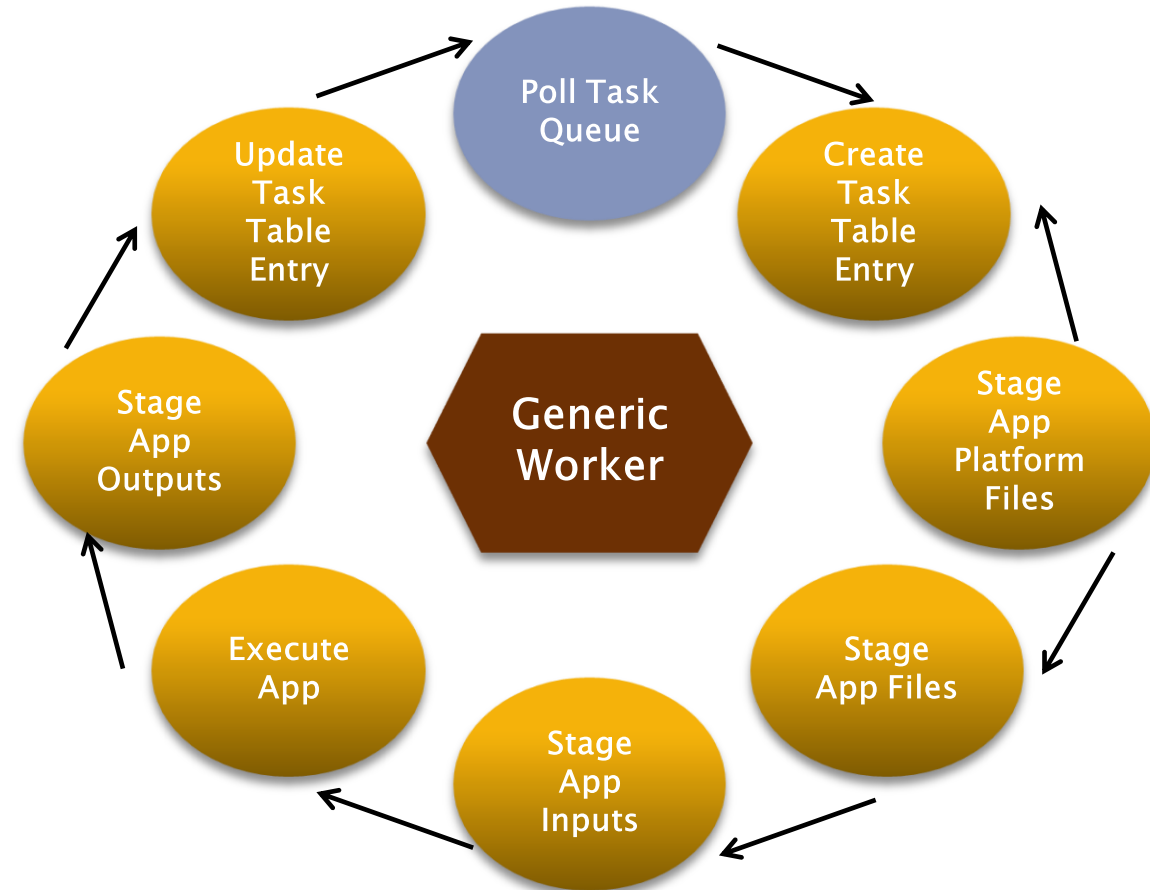


Reduction Service (Only One Stage Shown)



Inside A Generic Worker

- ▶ Manages application sandbox
 - Ensures all application binaries such as the MatLab runtime are installed for “known” application types
 - Stages all input blobs from Azure storage to local files
 - Passes any marshalled inputs to uploaded application binary
 - Stages all output blobs to Azure storage from local files
 - Preserves any marshalled outputs to the appropriate Task table
- ▶ Manages all task status
 - Dequeues tasks created by the Service Monitor
 - Retries failed tasks 3 times
 - Maintains all task status
- ▶ Simplifies desktop development and cloud deployment



Pipeline Stage Interactions

- ▶ The Web Portal Role, Service Monitor Role and 5 Generic Worker Roles are deployed at most times
 - 5 Generic Workers are sufficient for reduction algorithm testing and development (\$20/day)
 - Early results returned to scientist while deploying up to 93 additional Generic Workers; such a deployment typically takes 45 minutes
 - Deployment taken down when long periods of idle time are known
 - Heuristic for scaling number of Generic Workers up and down
- ▶ Download stage runs in the deep background in all deployed generic worker roles
 - IO, not CPU bound so no competition
- ▶ Reduction tasks that have available inputs run preferentially to Reprojection tasks
 - Expedites interactive science result generation
 - If no available inputs and a backlog of reprojection tasks, number of Generic Workers scale up naturally until backlog addressed and reduction can continue
 - Second stage reduction runs only after all first stage reductions have completed
- ▶ Reduction results can be downloaded following emailed link to zip file



Download

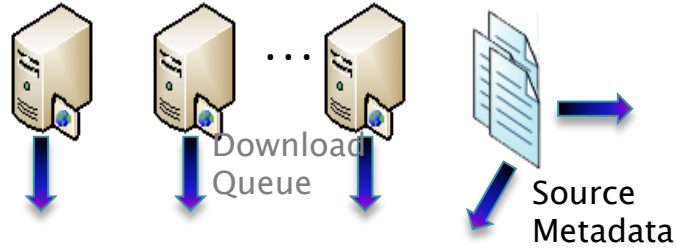
Reprojection

Reduction

Costs for 1 US Year ET Computation

- ▶ Computational costs driven by data scale and need to run reduction multiple times
- ▶ Storage costs driven by data scale and 6 month project duration
- ▶ Small with respect to the people costs even at graduate student rates !

Source Imagery Download Sites



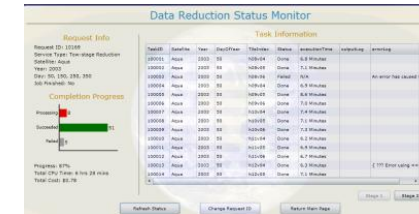
MODIS Data Reduction Service web portal interface showing fields for Year, Days, Satellites, Tiles, and Parameters, along with buttons for Submit Request and Return Main Page.

Request Queue



Scientists

AzureMODIS Service Web Role Portal



Scientific Results Download

Data Collection Stage

- 400–500 GB
- 60K files
- \$50 upload
- \$450 storage
- 10 MB/sec
- 11 hours
- <10 workers

Reprojection Queue

Reprojection Stage

- 400 GB
- 45K files
- \$420 cpu
- \$60 download
- 3500 hours
- 20–100 workers

Reduction #1 Queue

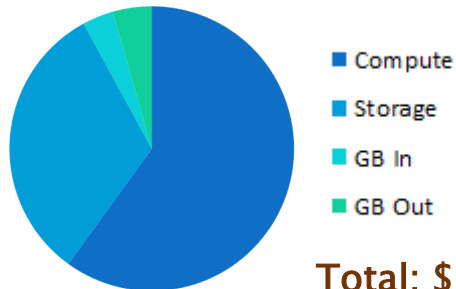
Derivation Reduction Stage

- 5–7 GB
- 5.5K files
- \$216 cpu
- \$1 download
- \$6 storage
- 1800 hours
- 20–100 workers

Reduction #2 Queue

Analysis Reduction Stage

- <10 GB
- ~1K files
- \$216 cpu
- \$2 download
- \$9 storage
- 1800 hours
- 20–100 workers



Total: \$1420

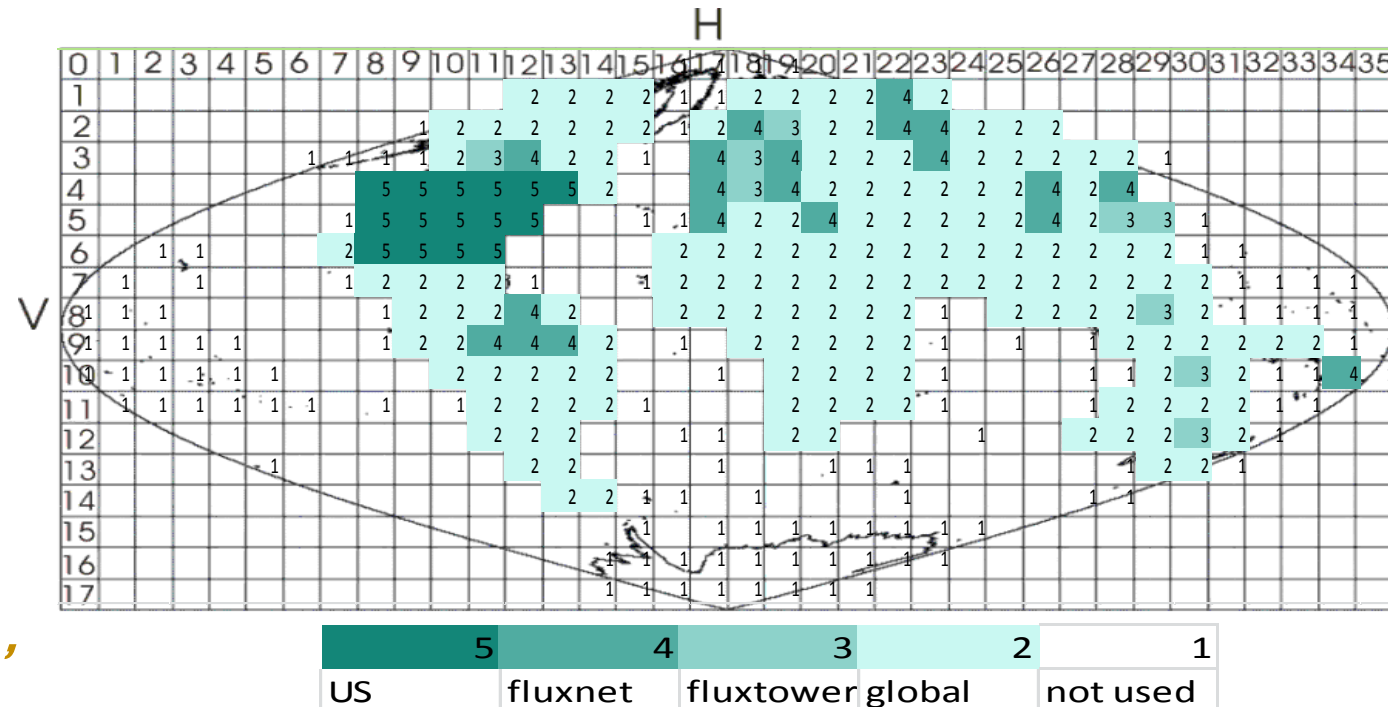
Scaling up to Global

*It appears to be monumental only because it's art.
Christo*

Sizing the 3 year MODIS Azure Global Computation

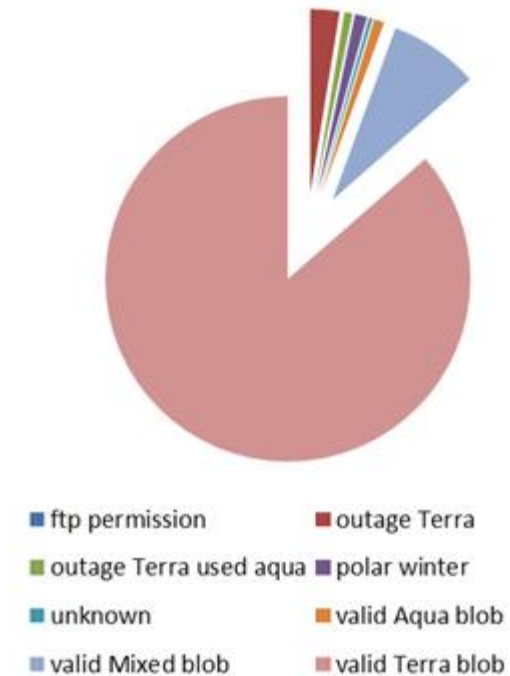
- ▶ 194 sinusoidal cells, each covers 1.2x1.2 KM or 11M 5 KM pixels)
- ▶ 1.06 M reprojected tiles and 40.5K source sinusoidal tiles
- ▶ 8 TB (>10 M files) downloaded from NASA ftp
- ▶ Not all files are downloaded or reprojected at first (3 rapid retries) attempt or actually available due to satellite outage, polar winter, missing tiles, etc etc.
- ▶ 55 NASA download days
- ▶ 150K reprojection compute hours
- ▶ 940 TB moved across Azure fabric
- ▶ 10 download result days (est) via IN2 bridge

*15 seconds on the Cray Jaguar (1.75 PFLOPs),
but only if we could get the PB in*



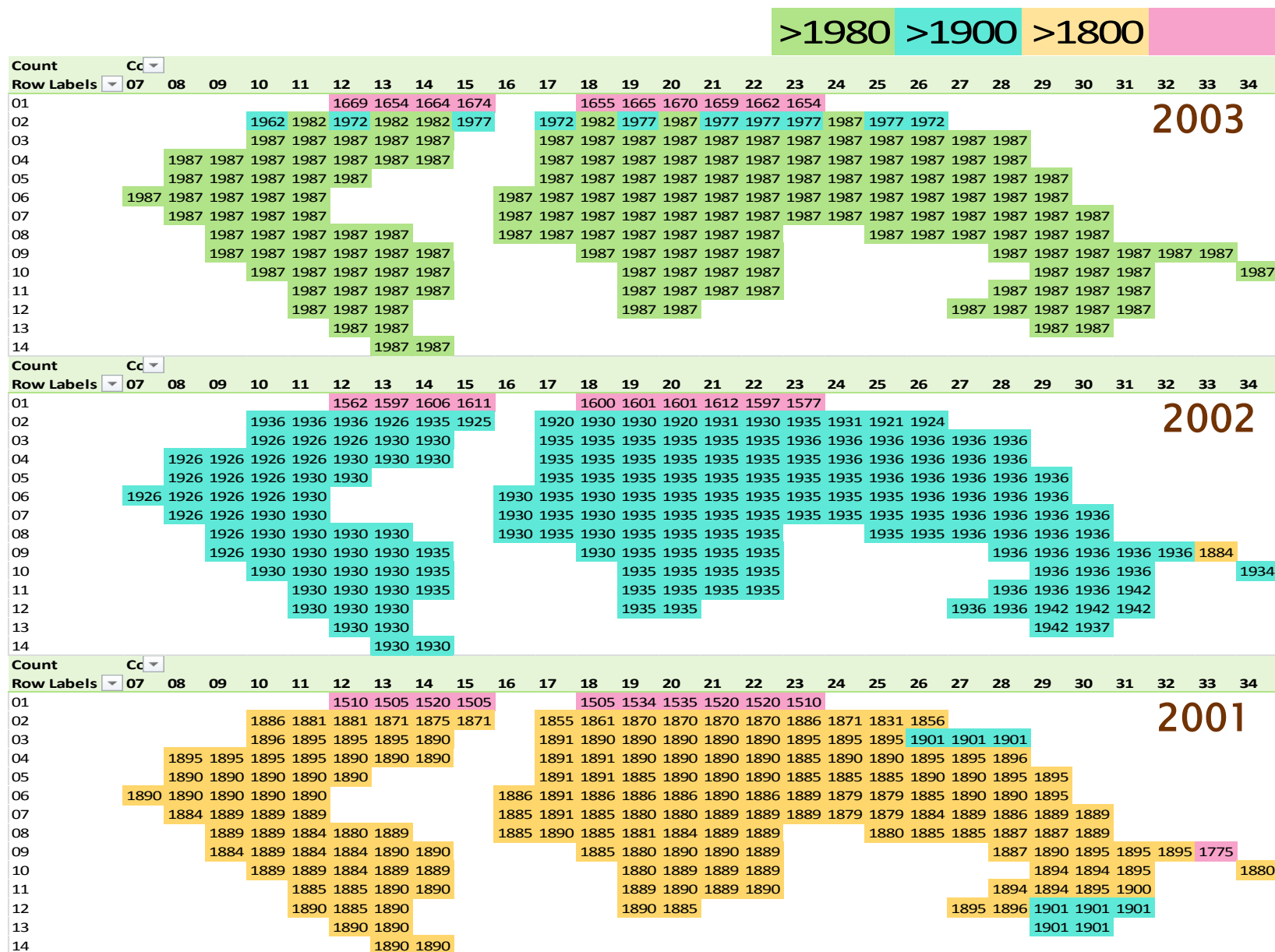
Input Tile File Tag Summary (1.1 M tiles)

- ▶ **Valid blob.** The tile is present in the blob store so has been downloaded or reprojected and will be used in the computation.
- ▶ **Download.** The sinusoidal projection tile is not in the blob store but could be downloaded from the NASA ftp site.
- ▶ **Reproject.** The reprojected tile is not in the blob store, but the swath precursors can be downloaded from the NASA ftp site and the a reprojection job run. This includes not only Terra tiles but also Aqua tiles that can be used to fill gaps in the Terra tiles.
- ▶ **Extra Aqua.** The reprojected valid blob tile is unnecessary for the computation as there is an existing valid or reprojectable Terra tile.
- ▶ **No Aqua.** The tile is not available because Aqua data is not available (eg MYD04 in 2001 or MCD15A2 in 2001).
- ▶ **No product.** The tile is not available because the product was not made (for reasons that usually, but not always includes pre-Aqua launch). Example is MOD44 in 2001.
- ▶ **Outage Terra.** The tile is not available due to a Terra satellite outage.
- ▶ **Outage Terra used Aqua.** The Terra tile is not available due to a Terra satellite outage ; an Aqua tile can be used.
- ▶ **Used Aqua.** The tile is a Terra tile for which we do not have a valid blob and cannot reproject to make one, but we can make an Aqua tile. The associated Aqua tile will either be tagged as “valid blob” (we have it) or “reproject “ (we need to make it).
- ▶ **Polar Winter.** The tile is in the v01 or v02 band and doy < 16 or doy > 339 and there is no swath precursor or sinusoidal projection tile on the NASA ftp site.
- ▶ **Unused Aqua.** The tile is an (to be reprojected) Aqua tile for which we do not have a valid blob, but we don’t need it because we have a valid Terra tile.
- ▶ **NASA protection.** The tile can be seen on the NASA ftp site, but all attempts to download fail with protection violation.
- ▶ **Unknown.** Tiles that are none of the above.



What's available ?

- ▶ Charts show all valid blob, download or reprojection tiles
- ▶ All tiles available => 1988 per year per cell
- ▶ Gaps most commonly due to satellite outages or polar winter
- ▶ Some transient gaps due to errors creating geo-spatial lookup or late addition of polar cells

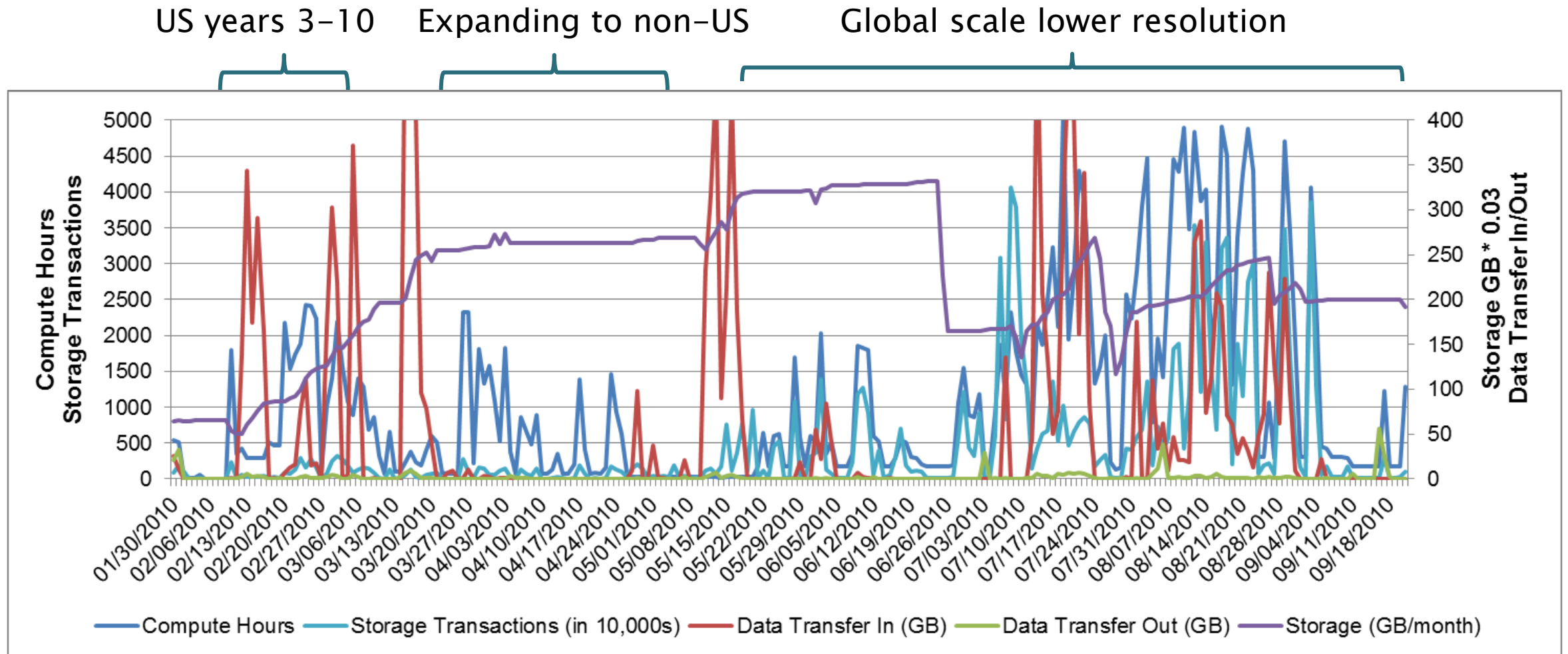


The “ity” Experience

*Even the stock holders of the phone company
hate the phone company!
Dr. Sidney Schaefer*

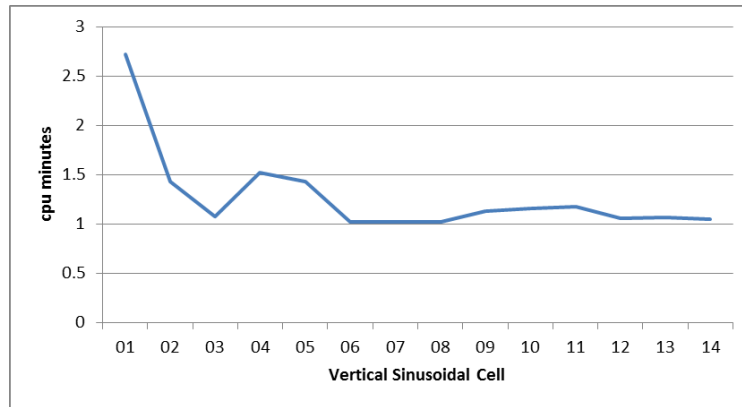
Agility

- ▶ The computation changed over time while Azure just scaled

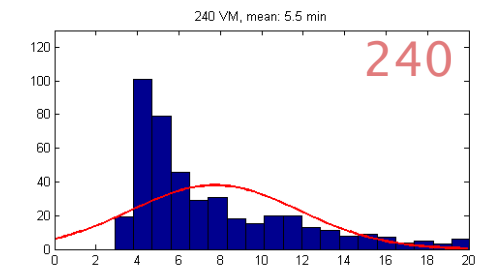
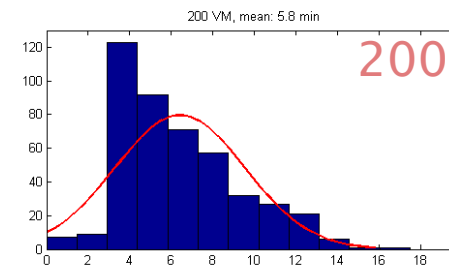
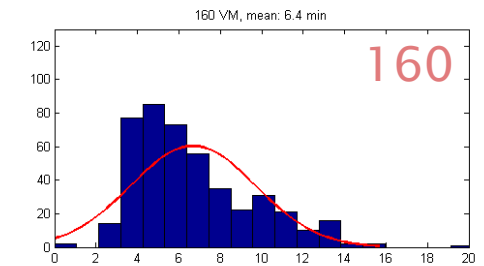
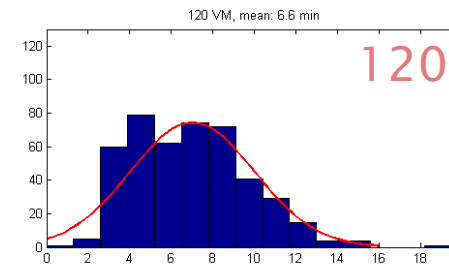
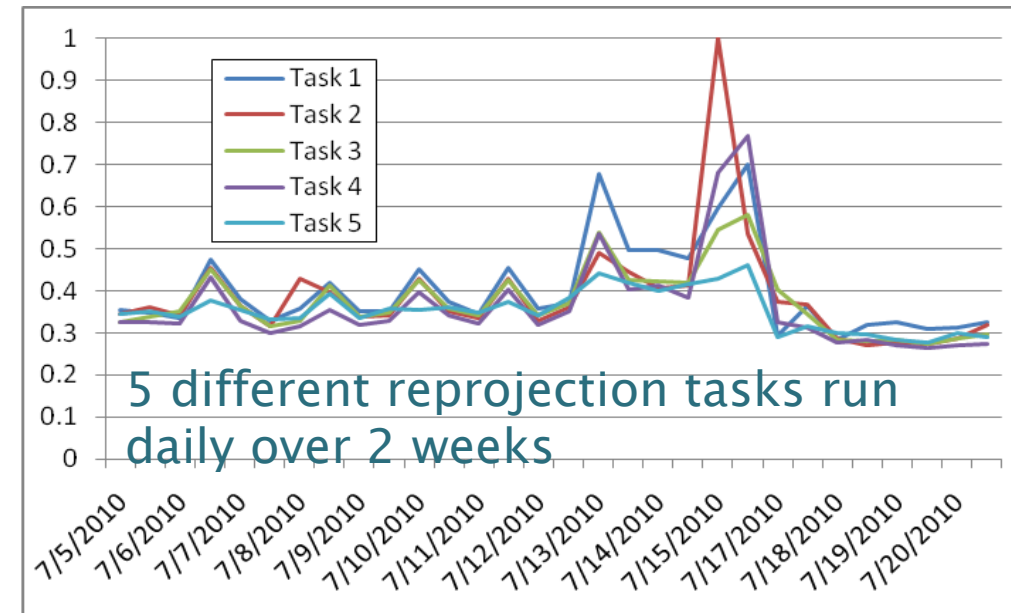


Predictability

- ▶ Performance varies over time: rerunning the same task gives different timings on different days
- ▶ Performance varies over space: satellites are over the poles more often



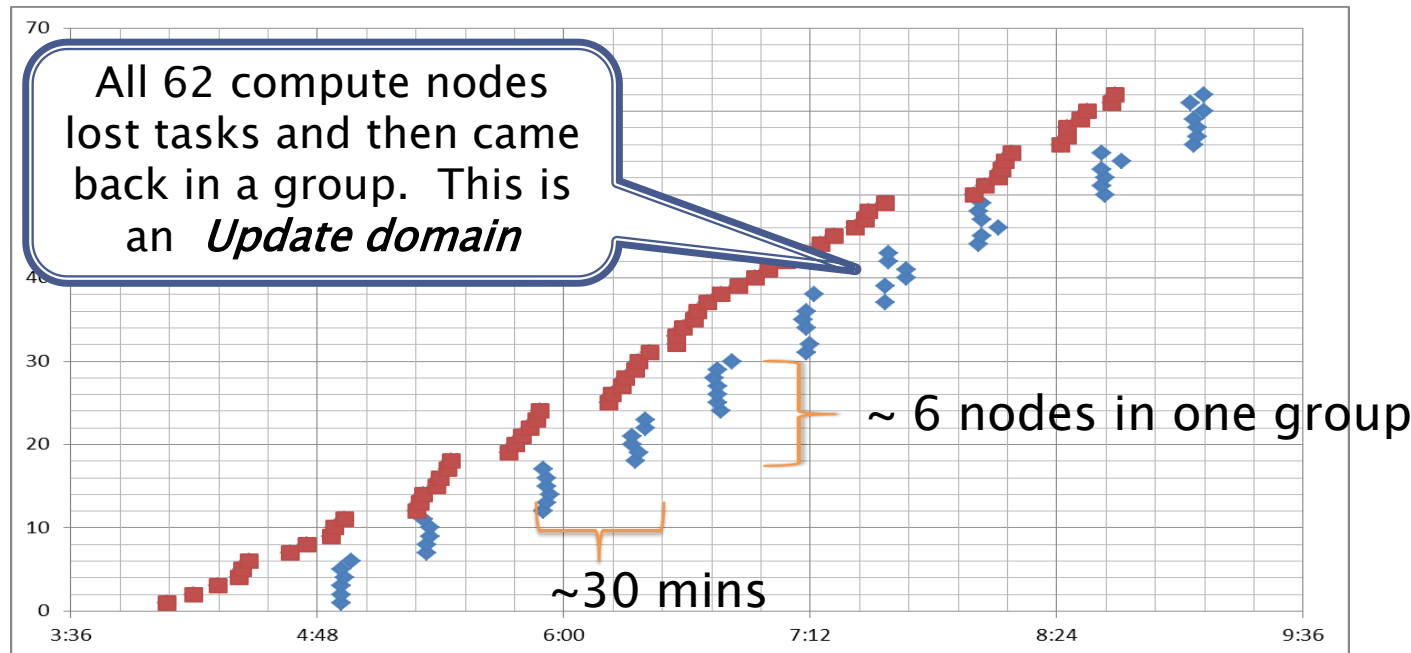
Average reprojection time (after algorithm improvements!) as a function of longitude



The same reduction task run on different numbers of VMs

Reliability

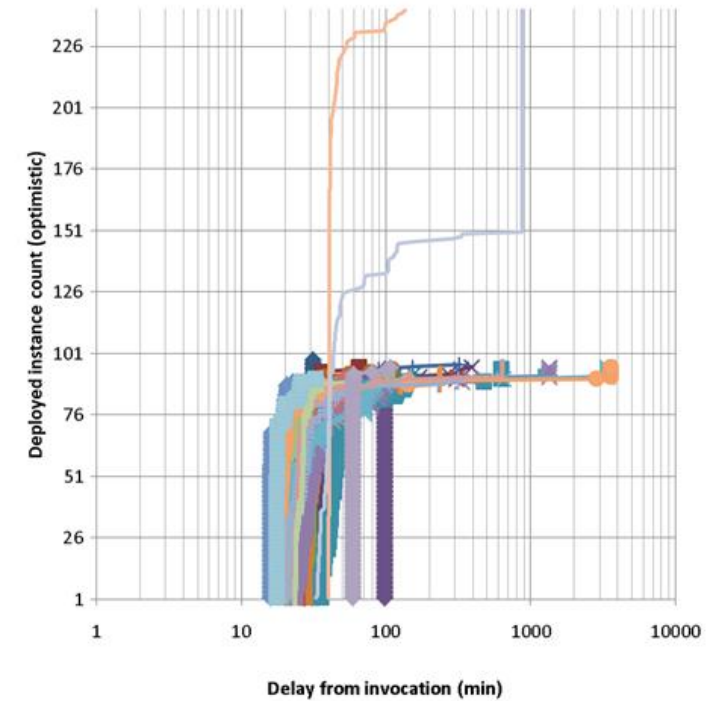
- ▶ Even with 99.999% reliability, bad things happen
 - 1–2 % of MODIS Azure tasks fail but succeed on retry



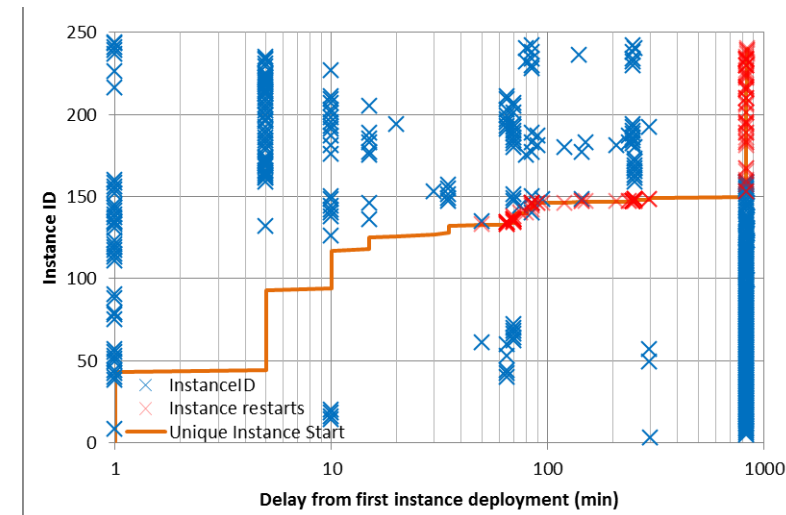
From AzureBlast

http://research.microsoft.com/en-us/people/barga/faculty_summit_2010.pdf

Observed VM starts for 76–100 VMs



Worst case attempt to start 250 VMs

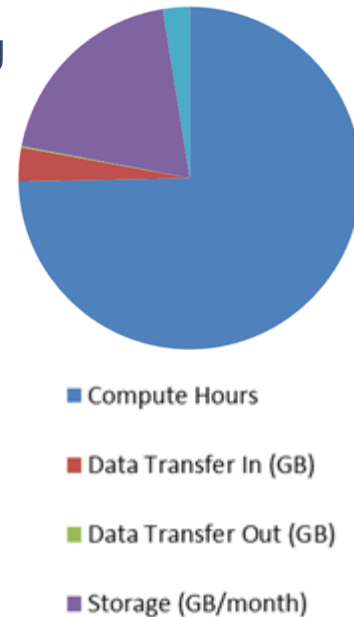


Fiscal Responsibility

- ▶ Billing is daunting
 - Neither we nor our academic collaborators are used to seeing bills
 - How *should* we think about them?
 - No billing cap means constant monitoring
- ▶ Billing is confusing
 - Instances are billed when deployed even if actually idle so comparing our usage log to the bill is at best approximate
 - Daily storage costs are amortized over the billing cycle so you must guestimate end cost
 - While you can ask for a refund, that takes a verified support call outage and time.
 - Online bill is autogenerated so must be accessed manually (no email)

Event Date	Name	Type	Region	Resource	Consumed	Sub Region	Service	Component	Service Info 1	Service Info 2	Additional Info
4/13/2010	Windows Azure Compute		North America	Compute Hours	73	South Central US	Compute	MODIS Data Services(Modis Data Service)			ComputeEmail
4/13/2010	Windows Azure Platform - All Services		North America	Data Transfer Out (GB)	0.001237	South Central US	Storage	MODIS Source Data Products			External
4/13/2010	Windows Azure Platform - All Services		North America	Data Transfer In (GB) - Off Peak	0.000001	South Central US	Storage	MODIS Source Data Products			External
4/13/2010	Windows Azure Platform - All Services		North America	Data Transfer Out (GB) - Off Peak	0.000018	South Central US	Storage	MODIS Source Data Products			External
4/13/2010	Windows Azure Platform - All Services		North America	Data Transfer Out (GB) - Off Peak	0.000044	South Central US	Storage	Reduction Results			External
4/13/2010	Windows Azure Platform - All Services		North America	Data Transfer Out (GB) - Off Peak	0.000002	South Central US	Storage	Reprojection Results			External
4/13/2010	Windows Azure Platform - All Services		North America	Data Transfer In (GB)	0.000032	South Central US	Compute	MODIS Data Services(Modis Data Service)			External
4/13/2010	Windows Azure Platform - All Services		North America	Data Transfer Out (GB)	0.000033	South Central US	Compute	MODIS Data Services(Modis Data Service)			External
4/13/2010	Windows Azure Platform - All Services		North America	Data Transfer In (GB) - Off Peak	0.000003	South Central US	Compute	MODIS Data Services(Modis Data Service)			External
4/13/2010	Windows Azure Platform - All Services		North America	Data Transfer Out (GB) - Off Peak	0.000003	South Central US	Compute	MODIS Data Services(Modis Data Service)			External
4/13/2010	Windows Azure Platform - All Services		North America	Data Transfer Out (GB)	0.000026	South Central US	Storage	Resources			External
4/13/2010	Windows Azure Platform - All Services		North America	Data Transfer In (GB)	0.000002	South Central US	Storage	MODIS Source Data Products			External
4/13/2010	Windows Azure Storage		North America	Storage (GB/month)	0.103332	South Central US	Storage	Resources			
4/13/2010	Windows Azure Storage		North America	Storage Transactions (in 10,000s)	0.2866	South Central US	Storage	Resources			
4/13/2010	Windows Azure Storage		North America	Storage (GB/month)	133.0917	South Central US	Storage	MODIS Source Data Products			
4/13/2010	Windows Azure Storage		North America	Storage Transactions (in 10,000s)	4.846	South Central US	Storage	MODIS Source Data Products			
4/13/2010	Windows Azure Storage		North America	Storage (GB/month)	14.84042	South Central US	Storage	Reduction Results			
4/13/2010	Windows Azure Storage		North America	Storage Transactions (in 10,000s)	0.0006	South Central US	Storage	Reduction Results			
4/13/2010	Windows Azure Storage		North America	Storage (GB/month)	92.30063	South Central US	Storage	Reprojection Results			
4/13/2010	Windows Azure Storage		North America	Storage Transactions (in 10,000s)	0.0006	South Central US	Storage	Reprojection Results			

One day of ModisAzure billing



100 instances @ \$0.12 per hour = \$288 per 24 hours

1 TB @ .15GB/mo = \$150.

Cumulative MODIS Azure billing (\$39K)

Maintainability

- ▶ Some “Early Adopter” artifacts
 - Generic worker sandbox
 - “dir” for blobs : need to have a parsable list, not just browse and many tools simply could not scale beyond $O(50K)$ blobs
 - “downloader” for blobs : early SDK utility retired by end of CTP.
 - Slow upload (FEDEX disk is still “in plan”; IN2 connections helped download tremendously
- ▶ Can we move catalog and other tracking to SQL Azure for better scaling?
 - Current tracking database is 140 GB
 - Partitions naturally, but would mean \$300/mo (external) charges.



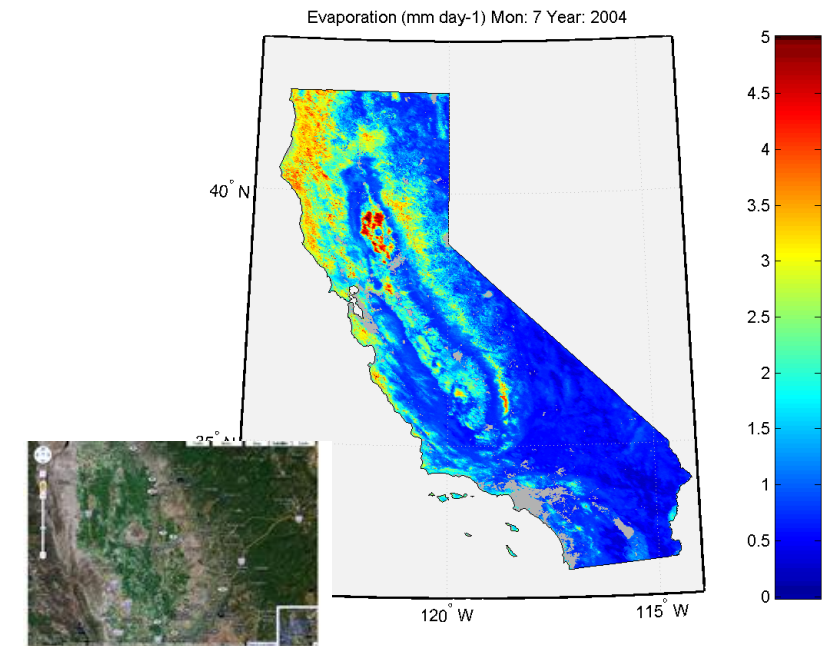
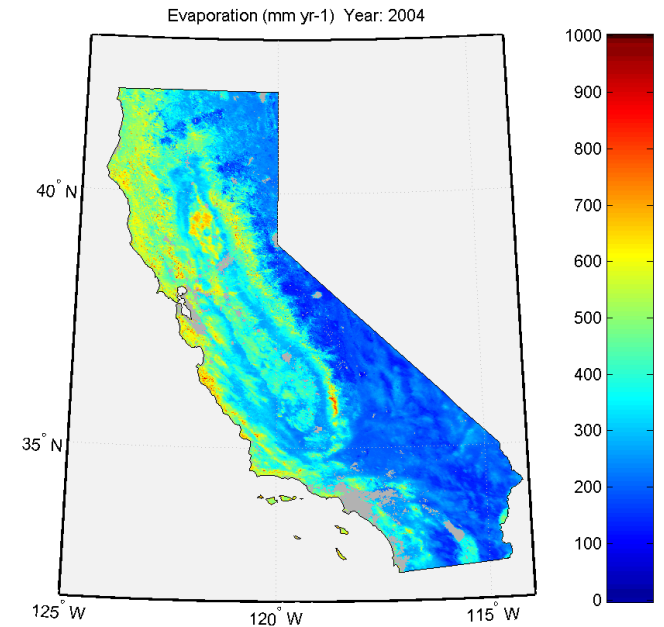
Summary and Prepublication Results

Baby, it's the beginning of a great adventure.

Lou Reed

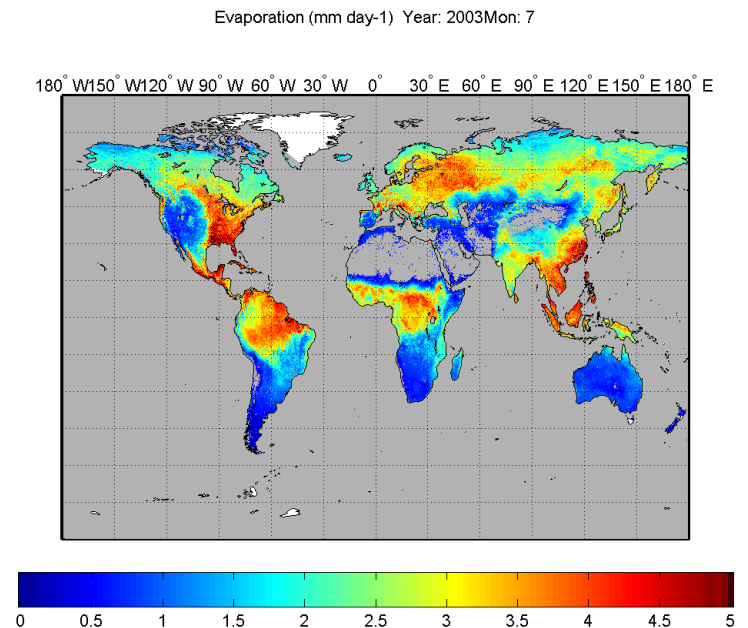
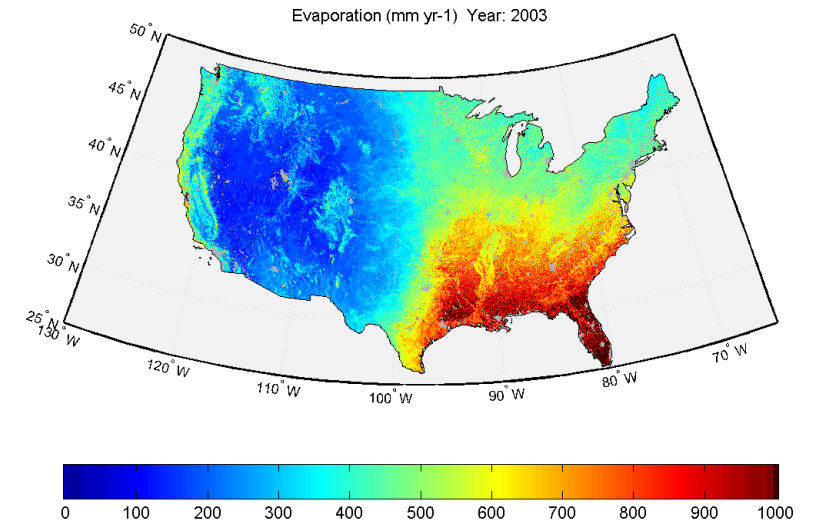
Azure Learnings

- ▶ Putting all your eggs in the cloud basket means watching that basket
 - Cloud scale resources often mean you still manage small numbers of resources: 100 instances over 24 hours = \$288 even if idle
 - Where is the long term archive for any results ?
- ▶ Azure is a rapidly moving target and unlike the Grid
 - Commercial cloud backed by large commercial development team
 - Current target applications are mid-range or smaller – MODIS Azure is currently at the fringe
- ▶ Scaling up requires additional work as understanding even a 0.01% failure rate is time consuming
 - Bake in the faults for scaling and resilience
 - Bake in the catalog for end:end reconciliation of sources and results



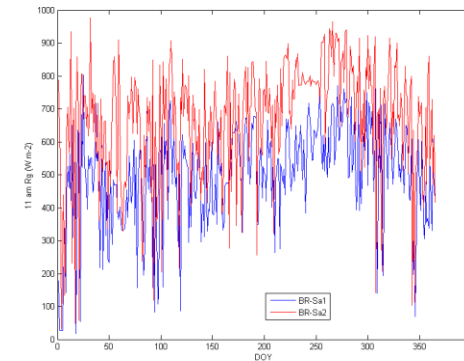
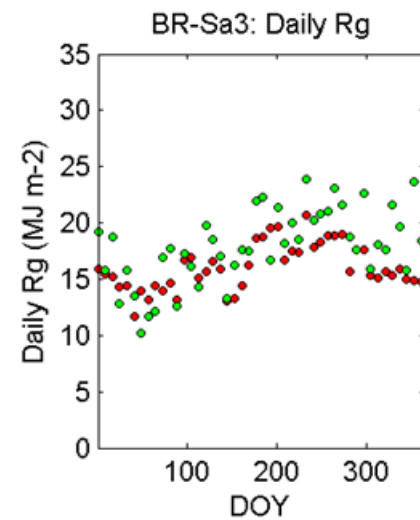
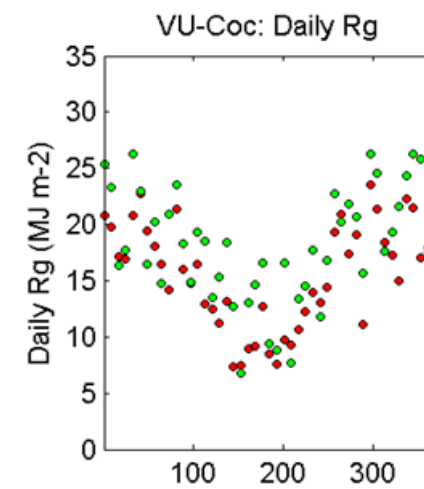
eScience Learnings

- ▶ Lowering the barriers to use remote sensing data can enable science
 - NASA makes the data accessible, not science ready
 - At AGU 2009, we learned that a cloud service that just made on-demand jpg mosaics would help tremendously
- ▶ Science and algorithm debugging benefit from the same infrastructure as both need to scale up and down
 - Debugging an algorithm on the desktop isn't enough – you have to debug in the cloud too
 - Whenever running at scale in the cloud, you must reduce down to the desktop to understand the results
- ▶ Scaling up means expanding the science
 - California, New England, and Florida are different
 - Boreal forests, savannahs, fertilization practices differ across the globe

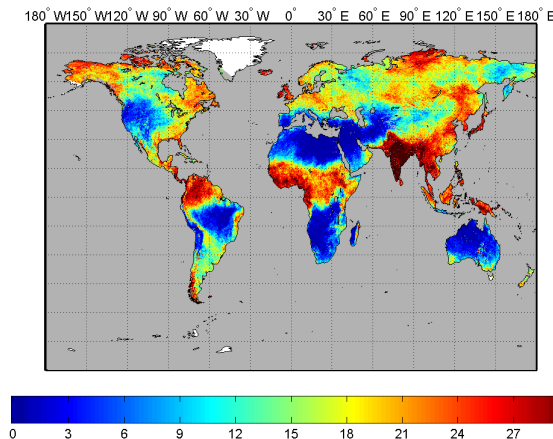


4th Paradigm Learnings

- ▶ Developing concrete plans for validation, sensitivity analysis, and mining prior to having results in hand is tricky
 - Precedents break down when scaling 100x or more
 - Sub-discipline familiarity a good start – our initial plans centered on FLUXNET tower data
 - Large sanity check aggregates a good start – our carbon fixation is in range of the literature estimates
 - Watershed aggregate comparison in the US crossed disciplines, length, and time scales as well as introducing yet different grunge.
- ▶ “Everybody knows” local knowledge plays a big role
 - Citizen science opportunity is anecdotal rather than quantified voting
 - Machine learning seems possible, but likely involves categorical geospatial subdivisions and some science cross checking



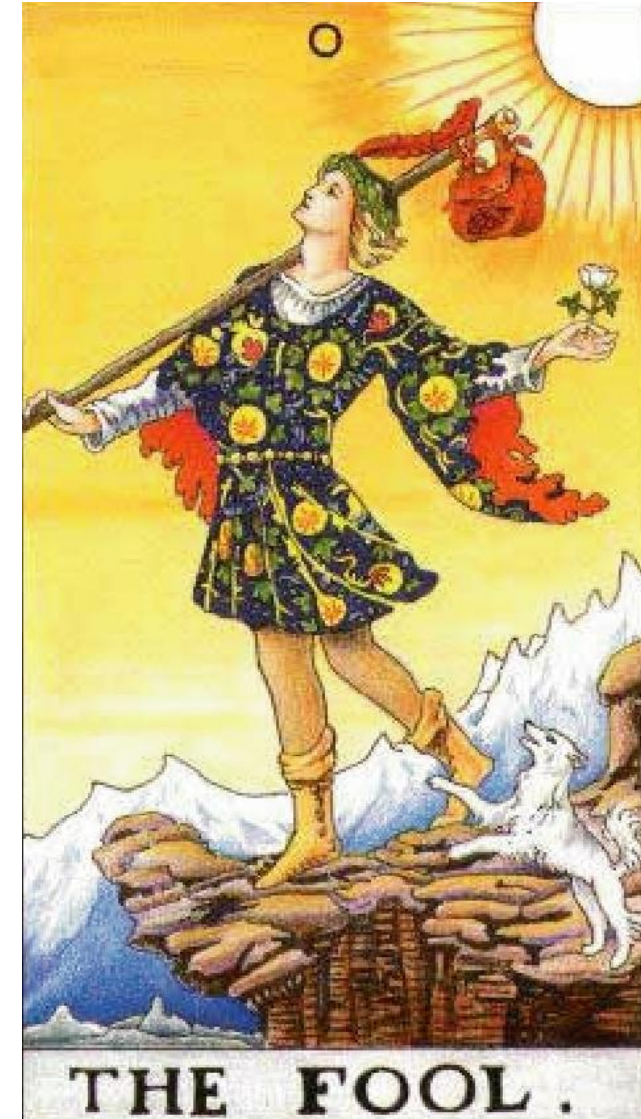
Cloudy days (days mon-1) Year: 2003Mon: 7



Current computation underestimates solar radiation in the tropics. Some of the comparison sites had to be discarded due to systematic drift of sensors used. And that data set has been used by others and passed the FLUXNET QA/QX.

4th Paradigm Challenges Forward

- ▶ How should we proceed to understanding the our global computation and related other computations well enough to improve such a computation over the next few years ?
 - Are there aggregate approaches such as computing statistics rather than values then statistics to reduce the overall computation requirements?
 - What should we do about the science factors we omitted such as elevation changes ?
 - What is the role for machine learning ? How can we engage?
- ▶ Since the dominant cost is people, how can we generalize the compute infrastructure to a wider class of computations?
 - Would Dryad/HPC/LINQ be faster, easier, more maintainable ?



Acknowledgements

- ▶ Scientists
 - Youngryel Ryu
 - Thomas Moran
 - Dennis Baldocchi
 - James Hunt
- ▶ Computer Scientists
 - Jie Li
 - You-Wei Cheah
 - Keith Jackson
 - Marty Humphrey
 - Deb Agarwal
 - Keith Beattie
- ▶ Others
 - The FLUXNET Collaboration (<http://www.fluxdata.org>)
 - Roger Barga
 - Dan Fay
 - Dennis Gannon
 - David Heckerman
 - Tony Hey
 - Yogesh Simmhan

Youngryel was lonely with 1 PC

